

On Particle Learning*

NICOLAS CHOPIN¹, ALESSANDRA IACOBUCCI², JEAN-MICHEL MARIN^{1,3},

KERRIE L. MENGENSEN⁴, CHRISTIAN P. ROBERT^{1,2,4},

ROBIN RYDER^{1,2}, AND CHRISTIAN SCHÄFER^{1,2}

¹CREST, Paris ²Université Paris-Dauphine, CEREMADE

³IM3, Université Montpellier 2 ⁴Queensland University of Technology

November 22, 2010

Abstract

This document is the aggregation of several discussions of Lopes et al. (2010) we submitted to the proceedings of the Ninth Valencia Meeting, held in Benidorm, Spain, on June 3–8, 2010, in conjunction with Hedibert Lopes’ talk at this meeting. The main point in those discussions is the potential for degeneracy in the particle learning methodology, related with the exponential forgetting of the past simulations. We illustrate the resulting difficulties in the case of mixtures.

Keywords: Attrition; degeneracy; evidence; importance sampling; Marginal likelihood; Markov chain Monte Carlo; mixtures of distributions; particle filter; sequential sampling; simulation.

1 The case of mixtures (Mengersen, Iacobucci and Robert)

In this discussion, we primarily consider the performances of the particle learning (PL) technique of Lopes et al. (2010) in the specific case of mixtures of distributions.

1.1 Particle learning

Reminiscing similar remarks made during Professor Polson’s talk at the ISBA 2008 World meeting on Hamilton Island, we do not understand the purpose of the dismissal of MCMC methods found in the paper (“more for less”, “direct approximations”, &tc.) Convergence of MCMC methods has been the core activity of many top researchers in the past two decades, first and foremost Gareth Roberts and Jeff Rosenthal, whose work cannot be so casually ignored! Especially when considering that, first, the main appeal in using particle methods (Gordon et al., 1993) is in handling massive data flux at frequencies MCMC cannot face—and this stands quite separate from a convergence issue—and, secondly, the body of work produced by the authors as listed in the reference list does not include any in-depth study of the convergence properties of the PL method.

As also argued in other discussions therein, the lack of warning in Lopes et al. (2010) or in previous papers about the unavoidable degeneracy of the method is more than puzzling, as the authors undoubtedly are aware of this. The short paragraph about Monte Carlo errors contained in the current paper can be construed to be misleading in this regard since the Monte Carlo error C_t/\sqrt{n} does not account for the resampling step. As demonstrated in the discussion by Robert and Ryder, the error may end up being $O(1/t)$ and miss the standard \sqrt{n} Monte Carlo convergence. The corpus of work thus produced so far seems to limit itself to the PL processing of an increasing sequence of state-space and dynamic

*N. Chopin and C.P. Robert are partly supported by the 2007–2010 grant ANR-07-BLAN-0237-01 “SP Bayes”. J.-M. Marin and C.P. Robert are partly supported by the 2009–2012 grant ANR-09-BLAN-0218 “Big’MC”. Robin Ryder is funded by a postdoctoral fellowship from the Fondation des Sciences Mathématiques de Paris. Christian Schäfer is supported by a PhD grant from CREST.

examples where the PL method does produce a reasonable output, but this series of case-studies does not constitute a sufficient validation in our eyes. (Some of the examples processed in the current paper are missing the hyperparameters chosen for their satisfactory resolution.)

We note as a side remark that the argument found in Section 1.5 of the current paper about the difficulty about the improper prior $p(\theta)$ being solved by using the mixture representation

$$p(\theta) = \int p(\theta|Z_0) p(Z_0) dZ_0$$

is nonsensical as currently stated: if the marginal distribution of θ is improper, so is the joint distribution. We also fail to understand where in the paper the authors manage to take a “new look at Bayes’s theorem”. If by this they mean the decomposition used in the first page of Section 1, this is a standard hidden Markov model property (Cappé et al., 2004).

1.2 Particle learning on mixtures

In the case of a mixture of k Poisson distributions,

$$f(x|\omega, \mu) = \sum_{i=1}^k p_i g(x|\lambda_i),$$

taken as a first example in Carvalho et al. (2009), the integrated predictive $p(y_{t+1}|\mathfrak{Z}_t)$ can be obtained in closed form, as detailed below (since this derivation is central to our own Monte Carlo experiment). For Poisson mixtures, the “essential” auxiliary variable is $\mathfrak{Z}_t = (n_1^t, \dots, n_k^t, s_1^t, \dots, s_k^t)$, where n_i^t denotes the number of observations allocated to the i -th component and s_i^t the sum of the observations allocated to component i ($1 \leq i \leq k$).

Thus, under a conjugate Dirichlet-Gamma prior assumption, using the delta function $\delta_{i=j}$ and the notation $\gamma_\cdot = \gamma_1 + \dots + \gamma_k$,

$$\begin{aligned} p(y_{t+1}|\mathfrak{Z}_t) &= \int p(y_{t+1}|\theta) p(\theta|\mathfrak{Z}_t) d\theta \\ &= \int \left\{ \sum_{i=1}^k p_i \frac{\lambda_i^{y_{t+1}} e^{-\lambda_i}}{y_{t+1}!} \right\} \mathcal{D}(p|(\gamma_j + n_j^t)) \prod_{j=1}^k \mathcal{G}(\lambda_j|\alpha_j + s_j^t, \beta_j + n_j^t) dp d\lambda \\ &= \int \sum_{i=1}^k \frac{\Gamma(\gamma_\cdot + t)}{\Gamma(\gamma_i + n_i^t)} \frac{\Gamma(\gamma_i + n_i^t + 1)}{y - t + 1! \Gamma(\gamma_\cdot + t + 1)} \frac{\Gamma(\alpha_i + s_i^t + y_{t+1})}{\Gamma(\alpha_i + s_i^t)} \frac{(\beta_i + n_i^t)^{\alpha_i + s_i^t}}{(\beta_i + n_i^t + 1)^{\alpha_i + s_i^t + y_{t+1}}} \\ &\quad \times \mathcal{D}(p|(\gamma_j + n_j^t + \delta_{i=j})) \prod_{j=1}^k \mathcal{G}(\lambda_j|\alpha_j + s_j^t + \delta_{i=j} y_{t+1}, \beta_j + n_j^t + \delta_{i=j}) dp d\lambda \\ &= \sum_{i=1}^k \frac{\Gamma(\gamma_\cdot + t)}{\Gamma(\gamma_i + n_i^t)} \frac{\Gamma(\gamma_i + n_i^t + 1)}{y - t + 1! \Gamma(\gamma_\cdot + t + 1)} \frac{\Gamma(\alpha_i + s_i^t + y_{t+1})}{\Gamma(\alpha_i + s_i^t)} \frac{(\beta_i + n_i^t)^{\alpha_i + s_i^t}}{(\beta_i + n_i^t + 1)^{\alpha_i + s_i^t + y_{t+1}}} \\ &= \sum_{i=1}^k \frac{\gamma_i + n_i^t}{\gamma_\cdot + t + 1} \frac{1}{y_{t+1}!} \frac{\Gamma(\alpha_i + s_i^t + y_{t+1})}{\Gamma(\alpha_i + s_i^t)} \frac{(\beta_i + n_i^t)^{\alpha_i + s_i^t}}{(\beta_i + n_i^t + 1)^{\alpha_i + s_i^t + y_{t+1}}} \end{aligned}$$

gives a closed form of the predictive distribution of y_{t+1} given the sufficient (simulated) statistic \mathfrak{Z}_t .

Furthermore, the distribution of the first version of the auxiliary variable, z_1 , is easily derived, as

$$\begin{aligned} \mathbb{P}(z_1 = j|y_1) &\propto \int p_j g(y_1|\lambda_j) \pi(p, \lambda) dp d\lambda \\ &\propto \int \lambda_j^{y_1 + \alpha_j - 1} e^{-\lambda_j(1 + \beta_j)} p_j^{1 + \gamma_j - 1} (1 - p_j)^{1 + \gamma_\cdot - \gamma_j - 1} dp_j d\lambda_j \\ &= \Gamma(y_1 + \alpha_j) (1 + \beta_j)^{-(y_1 + \alpha_j)} \Gamma(1 + \gamma_j) \Gamma(1 + \gamma_\cdot - \gamma_j) / \Gamma(2 + \gamma_\cdot) \end{aligned}$$

while the next allocations can be found as

$$\begin{aligned}\mathbb{P}(z_{t+1} = j | y_{t+1}, \mathfrak{Z}_t) &\propto p(y_{t+1}, k_{t+1} = j | \mathfrak{Z}_t) = \int p(y_{t+1}, k_{t+1} = j | \theta) \pi(\theta | \mathfrak{Z}_t) d\theta \\ &\propto \frac{(\gamma_j + n_j^t) \Gamma(\alpha_i + s_i^t + y_{t+1})}{\Gamma(\alpha_i + s_i^t)} \frac{(\beta_i + n_i^t)^{\alpha_i + s_i^t}}{(\beta_i + n_i^t + 1)^{\alpha_i + s_i^t + y_{t+1}}}\end{aligned}$$

and simulating the parameters of the Poisson model given \mathfrak{Z}_t is obvious, provided one uses a conjugate prior (Diebolt and Robert, 1990).

When applying both PL and MCMC techniques to a sample of size 10^4 extracted from the Monte Carlo study detailed in the discussion by Iacobucci et al., we obtain the output represented in Figure 1 for the posterior distributions of the λ_i 's. (Both PL and MCMC samples were re-ordered in terms of the λ_i 's. The plots are therefore the posterior distributions of the order statistics $\lambda_{(i)}$.) The discrepancy between both approaches is clear on this example. (We stress that this represents a “worst case” in the sense that the observations were chosen from the above-mentioned Monte Carlo experiment by selecting the sample producing the largest discrepancy in the evidence approximations. Random picks of samples from the Poisson mixture usually produce a better agreement.)

2 On the approximation of evidence (Iacobucci, Robert, Marin and Mengersen)

In this discussion, we consider the performances of the particle learning (PL) technique in the specific setting of mixtures of distributions and for the approximation of the *evidence*

$$\mathfrak{Z}_i = \int_{\Theta_i} \pi_i(\theta_i) f_i(y | \theta_i) d\theta_i,$$

aka the marginal likelihood. Through a simulation experiment, we examine how much the degeneracy that is inherent to particle systems impacts this approximation (We refer the reader to Chen et al., 2000, for a general approach to the approximation of evidence and to both Chopin and Robert, 2010, and Marin and Robert, 2010, for illustrations in the particular setting of mixtures.)

2.1 Approximation of the evidence

In the case of a mixture of k Poisson distributions,

$$f(x | \omega, \mu) = \sum_{i=1}^k p_i g(x | \lambda_i),$$

taken as an example in Lopes et al. (2010), and studied in Carvalho et al. (2009) the integrated predictive can be obtained in closed form, as derived in the discussion of Mengersen et al. This implies that the product approximation to the evidence

$$p(y^t) = \prod_{r=1}^t p(y_r | y^{r-1}) \approx \prod_{r=1}^t \frac{1}{N} \sum_{i=1}^N p(y_r | \mathfrak{Z}_{r-1}^{(i)})$$

proposed in Carvalho et al. (2009) and Lopes et al. (2010) can be implemented here. We thus use the setting of Poisson mixtures to evaluate this PL approximation of the evidence and we re-evaluate Carvalho et al.'s (2009) assessment that this “*approach offers a simple and robust sequential Monte Carlo alternative to the traditionally hard problem of approximating marginal predictive densities via MCMC output*”.

We note that, since the PL sample is considered as an approximate sample from the posterior $\pi(p, \lambda | y^t)$ it is possible to evaluate the evidence using Chib's (1995) formula rather than the above proposal of the authors. The availability of an alternative estimator of the evidence allows for a differentiation between

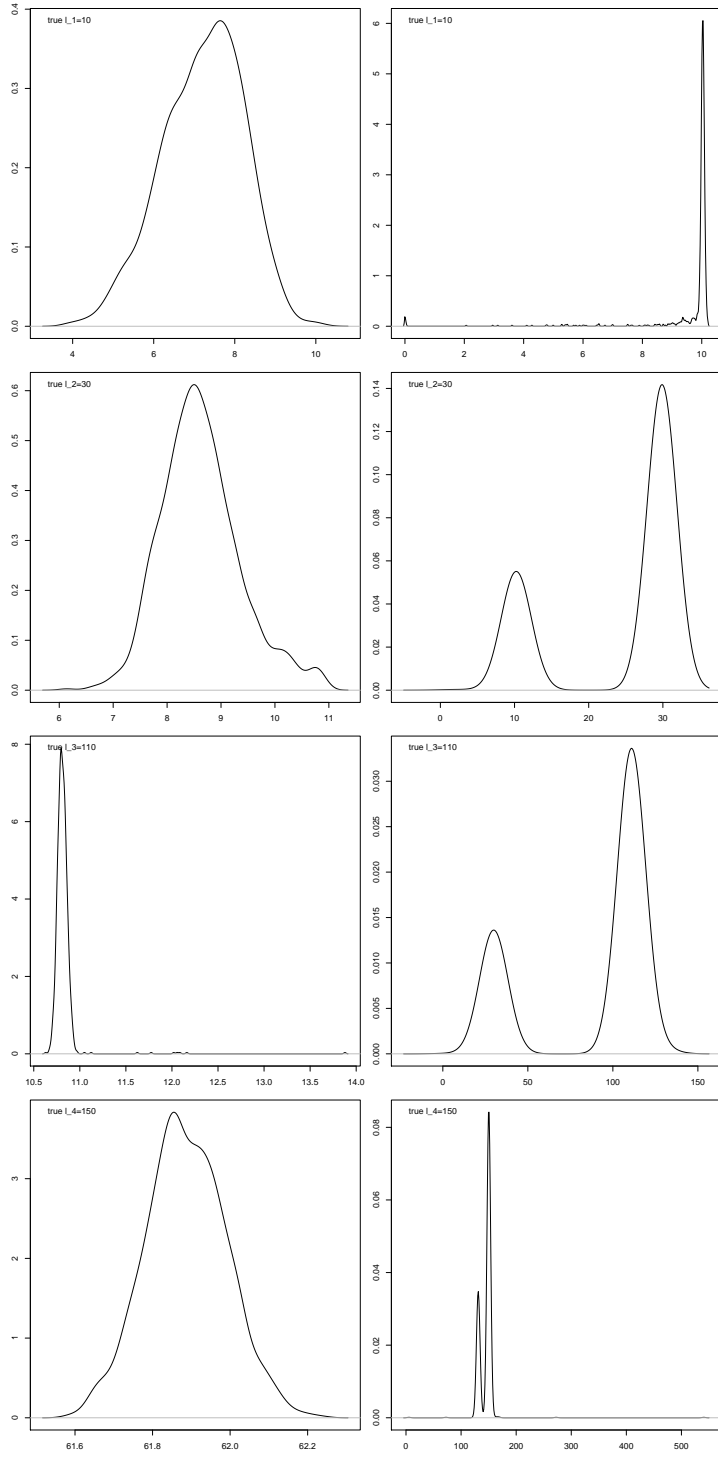


Figure 1: Comparison of the posterior distributions on the ordered λ_i 's produced by PL (*left*) and MCMC (*right*) for 10^4 simulated observations from a 4 component Poisson mixture and 10^4 particles/MCMC iterations. The curves are obtained by apply the R function `density()` to both samples. The true values are indicated at the top of each graph

the evaluation of approximation [of the target posterior distribution] resulting from the particle system (seen through a possible bias in Chib’s, 1995, version) and the evaluation of the approximation [of the evidence] resulting from the use of the product marginal in Lopes et al. (2010). Thus, in contrast to the other discussions of ours, we evaluate here the specific degeneracy of the evidence approximation due to using a product of approximations.

2.2 A Monte Carlo experimentation

In order to evaluate the performances of the PL algorithm when compared with the vanilla Gibbs sampler (Diebolt and Robert, 1990, 1994), we simulated 250 samples of size 10^4 from Poisson mixtures with 4 and 5 components and with either widely spaced or close components, $\lambda = (10, 50, 110, 150, 180, 210)$ and $\lambda = (10, 15, 20, 25, 30, 35)$, respectively, and with slightly decreasing weights p_i . We ran a 10^4 iteration Gibbs sampler for Figures 2–5, performing a further 10^6 iterations as a check of the stability of the MCMC approximation. (For Chib’s approximation to perform correctly, as noted in Berkhof et al., 2003 and Marin and Robert, 2010, it is necessary to average over all $k!$ permutations of the component indices for both the original PL sample and the MCMC sample in order to escape label switching issues.)

The first interesting outcome of our experiment is that the PL sample does not suffer from degeneracy for a small enough number of observations, since the ranges of the Chib’s (2005) approximations for both PL and MCMC samples (represented by the second and third columns in the boxplots) are then the same. However, as predicted by the theory (see the discussions by Chopin and Robert, and by Robert and Ryder), increasing the number of observations without simultaneously and exponentially increasing the number of particles necessarily leads to the degeneracy of the simulated sufficient statistic paths. In our experiment, this degeneracy always occurs between 5,000 and 10,000 observations. The phenomenon clearly appears on Figures 2–5 where both the range and the extremes of the evidence approximations significantly differ on the right hand side boxplot graph. (Again, the stability of the MCMC range was tested by running the Gibbs sampler for much longer and observing no variation.) This divergence is to be contrasted with Figure 1 in Carvalho et al. (2009) which concludes to an agreement between all approximations to the Bayes factor.

The second result that is relevant for our discussion is that the new approximation to the evidence proposed by the authors suffers from a severe bias as one proceeds through the observations. This issue is apparently unrelated to the degeneracy phenomenon observed above in that the discrepancy starts from the beginning, the closest approximation occurring for $n = 1,000$ observations. Note that Carvalho et al. (2009) mention that the evidence approximation based on particle learning was less variable. While this feature is not visible in our experiment, it is not necessarily a positive feature in any case, as shown in the current experiment. (In order to provide a better rendering of the comparison between the PL and the MCMC algorithms, we excluded the outliers from all boxplots. We however stress that both PL approaches had a higher propensity to outlying behaviour.) In the strongest case of discrepancy between PL and MCMC found in our experiment, Figure 6 illustrates the departure between the three approaches from a particularly influential observation, since the graphs are compared in terms of evidence *per observation*.

We thus conclude at the lack of robustness of the new approximation of evidence suggested in both Carvalho et al. (2009) and Lopes et al. (2010) (besides providing a reinforced demonstration of the overall difficulty with degeneracy).

3 Repeatability of the degeneracy (Iacobucci, Marin and Robert)

Following the floor discussion at the conference, we want to point out here that the divergence between the evidence evaluations observed in the discussion of Iacobucci et al. is not the result of an outlying Monte Carlo experiment but indeed a distributional property. This can be seen on Figure 7 which reproduces the study of Iacobucci et al. (in this set of comments) on the variation of the evidence in the specific setting of mixtures of Poisson distributions. For two given datasets, we repeated 683 times the three evidence approximations using the method proposed in Lopes et al. (2010), and Chib’s (1995) method applied to both the PL and MCMC samples. The divergence between the three evaluations is consistent across simulations, so repeating simulations does not help in exhibiting this divergence.

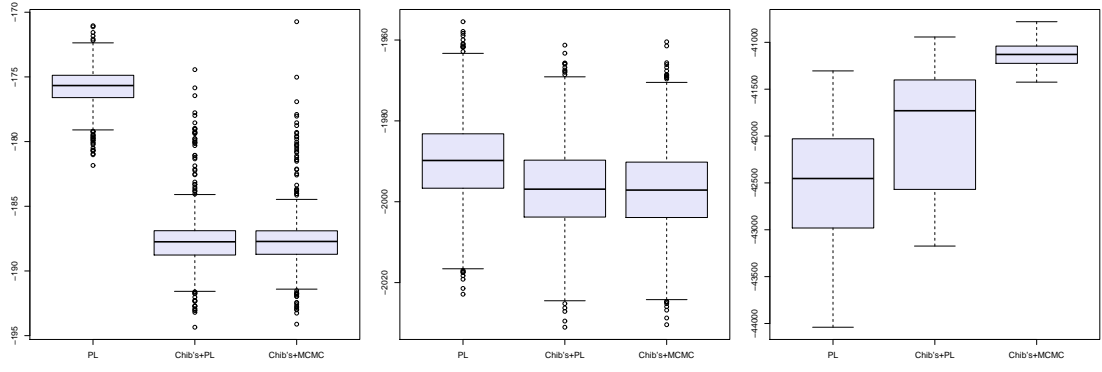


Figure 2: Evolution against the number of observations ($n = 100$, $n = 1,000$ and $n = 10,000$) of the evidence approximation based on a PL sample and Lopes et al. (2010) approximation, on a PL sample and Chib's (1995) approximation, on an MCMC sample and Chib's (1995) approximation, for a particle population of size 10,000, a mixture with 4 components and scale parameters $\lambda = (10, 50, 110, 150)$.

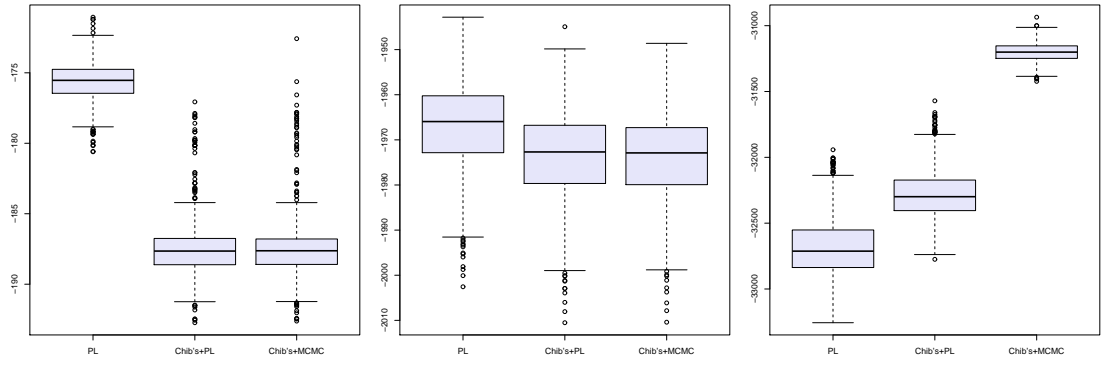


Figure 3: Same caption as Figure 2 for $\lambda = (10, 15, 20, 25)$.

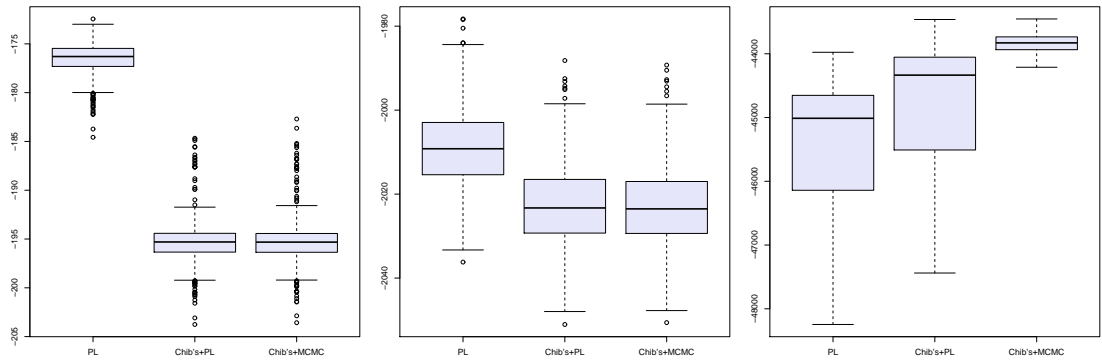


Figure 4: Same caption as Figure 2 for 5 components.

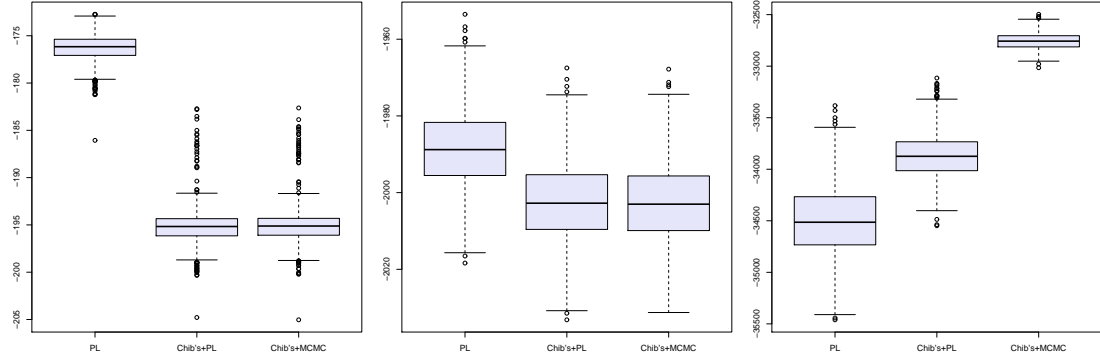


Figure 5: Same caption as Figure 3 for 5 components.

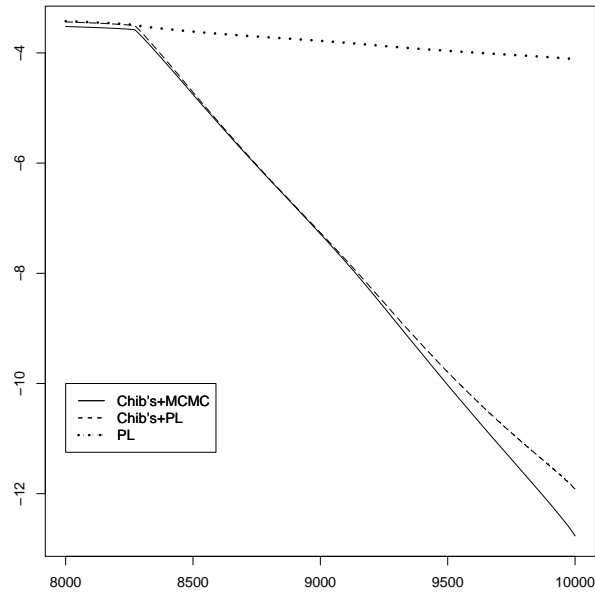


Figure 6: Evolution of the three approximations of the evidence *per observation* against the number of observations for a specific sample simulated from the same Poisson mixture as in Figure 2.

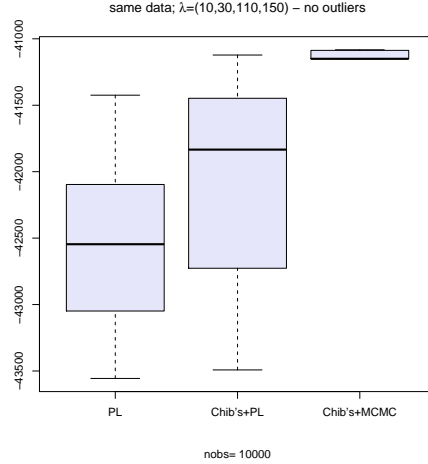


Figure 7: Range of the evidence approximation based on a PL sample and Lopes et al. (2010) approximation, on a PL sample and Chib’s (1995) approximation, on an MCMC sample and Chib’s (1995) approximation, for a particle population of size 10,000, a mixture with 4 components and scale parameters $\lambda = (10, 50, 110, 150)$, and 683 replications.

4 On degeneracy (Chopin and Robert)

In this discussion, we consider the performance of the particle learning technique of Lopes et al. (2010) in a limiting case, in order to illustrate the fact that a particle system cannot but degenerate, even when considering sufficient statistics Z_t with fixed dimensions.

4.1 Particle system degeneracy

When Lopes et al. (2010) state that $p(Z^t|y^t)$ is not of interest as the filtered, low dimensional $p(Z_t|y^t)$ is sufficient for inference at time t , they seem to implicitly imply that the restriction of the simulation focus to a low dimensional vector is a way to avoid the degeneracy inherent to all particle filters (see, e.g., Del Moral et al., 2006). However, the degeneracy of particle filters is an unavoidable consequence of the explosion of the state vector Z^t and the issue does not vanish because one is only interested in the marginal

$$p(Z_t|y^t) = \int p(Z^t|y^t) dZ^{-t}.$$

Indeed, as shown by the pseudo-code rendering in Lopes et al. (2010), the way PL produces a sample from $p(Z_t|y^t)$ is by sequentially simulating Z^t and by extracting Z_t as the final output from this sequence. The PL algorithm therefore relies on an approximation of $p(Z^t|y^t)$ and the fact that this approximation quickly degenerates as t increases, as discussed below and in the companion discussion by Robert and Ryder, obviously has an impact on the approximation of $p(Z_t|y^t)$.

Inherently, particle learning (PL) is at its core an auxiliary particle filter (Pitt and Shephard, 1999) applied in settings where there exists a sufficient statistic (Darmois, 1935) of reduced (or, even better, with fixed) dimension. The simulation scheme thus relies on resampling (Rubin, 1988, Kitagawa, 1996) for adjusting the distribution of the current particle population to the new observation y_{t+1} . Because of this continual resampling, the number of different values of Z_p ($p \geq 1$) contributing to the sufficient statistic Z_t ($t > p$) is decreasing in t at an exponential rate for a fixed p . Therefore, unless the size of the particle population exponentially increases with t (see Douc et al., 2002, and the companion discussion by Chopin and Schäfer), the sample of Z_t ’s will not be distributed as an iid sample from $p(Z_t|y^t)$. The following section very clearly makes this point through a simple if representative example.

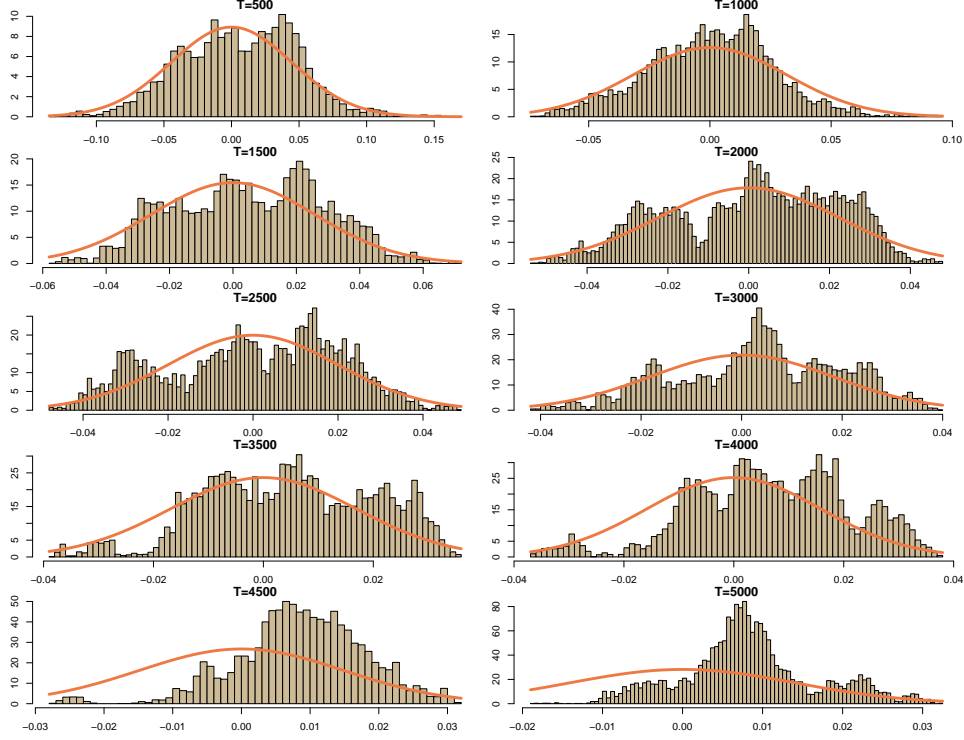


Figure 8: Evolution of the particle learning sample against the target distribution in terms of the number $T = 50, \dots, 5000$ of iterations, for a particle population of fixed size 10^4 .

4.2 A simple particle learning example

Consider the ultimate case where the z_t 's are completely independent from the observations y_t , $z_t \sim \mathcal{N}(0, 1)$, and where the empirical average of the z_t 's is the sufficient statistic. In this setting, the PL algorithm simplifies into the following iteration t :

1. Resample uniformly from (Z_1^t, \dots, Z_n^t) to produce $(\mathfrak{Z}_1^t, \dots, \mathfrak{Z}_n^t)$;
2. Generate $z_{it} \sim \mathcal{N}(0, 1)$;
3. Update $Z_i^{t+1} = (t\mathfrak{Z}_i^t + z_{it})/(t+1)$

The target distribution of the (sufficient) empirical average

$$Z^t = (z_1 + \dots + z_t)/t$$

is obviously the normal $\mathcal{N}(0, 1/t)$ distribution. A straightforward simulation of the above particle system shows how quickly the degeneracy occurs in the sample: Figures 8–9 show a complete lack of fit to the target distribution as early as $t = 500$ simulations when using 10,000 particles.

4.3 Conclusion

The paper Lopes et al. (2010) fails to mention the well-documented issue of particle degeneracy (Cappé et al., 2004, Del Moral et al., 2006), thus giving the impression that PL escapes this problem. Our simple example shows that a particle system cannot be expected to withstand an indeterminate increase in the number of observations without imposing a corresponding exponential increase in the particle size.

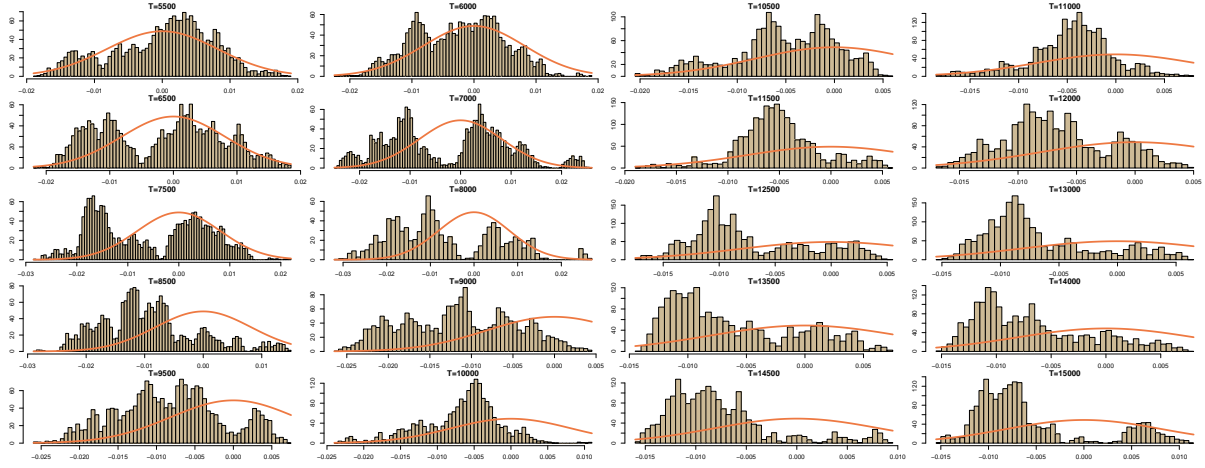


Figure 9: Figure 8 continued for $T = 5000, \dots, 15000$ iterations.

5 On the degeneracy of sufficient statistics (Robert, Ryder and Chopin)

In connection with the discussion of Chopin and Robert, we detail in this discussion how the degeneracy dynamics of the particle learning technique of Lopes et al. (2010) impacts the distribution of the sufficient (or “essential state vector”) statistics.

Lopes et al. (2010) focus on the distribution of a sufficient statistic, $p(Z_t|y^t)$, at time t . By insisting both on the low dimensionality of Z_t and on the sufficiency, they give the reader the impression that the poor approximation of the state vector Z^t resulting from the resampling propagation scheme does not impact $p(Z_t|y^t)$, since their statement “*at time T , PL provides the filtered distribution of the last essential state vector Z_T , namely $p(Z_T|y^T)$* ” (Section 1.2) does not mention any deterioration in the approximation—this is how we understand *filtered*—provided by PL. Because particle learning is inherently a particle filter (Pitt and Shephard, 1999), this intuition is unfortunately wrong, as shown below in the case of an empirical average of the past auxiliary variables Z_t . Contrary to the belief that “*resampling (...) is fundamental in avoiding a decay*” (Section 1.2), resampling necessarily leads to degeneracy unless the size of the particle population increases exponentially with t .

We thus consider again the case introduced by Chopin and Robert in their discussion, when the auxiliary variables $z_t \sim \mathcal{N}(0, 1)$ are independent from the observations y_t and where the essential state vector statistic is the empirical average of the z_t ’s. In this case, the distribution of the empirical average

$$Z^t = (z_1 + \dots + z_t)/t$$

is the normal $\mathcal{N}(0, 1/t)$ distribution, but the particle population degenerates into a single path from the point of view of this sufficient statistic. In other words, degeneracy occurs much faster than the root T forgetting of the past of the particle path that is due to the averaging. In order to support this perspective, we provide here a derivation of the variance of the particle population after t iterations.

Using the same notations as in Chopin and Robert, since $\mathbb{E}[Z_i^t] = 0$, $\text{var}(Z_i^t) = 1/t$ and $Z_i^t = \frac{t-1}{t}\mathfrak{Z}_i^t + \frac{z_{it}}{t}$, we consider

$$\begin{aligned} \mathbb{E}[Z_i^t Z_j^t] &= \left(\frac{t-1}{t}\right)^2 \mathbb{E}[\mathfrak{Z}_i^t \mathfrak{Z}_j^t] \\ &= \left(\frac{t-1}{t}\right)^2 (\mathbb{P}[\mathfrak{Z}_i^{t-1} = \mathfrak{Z}_j^{t-1}] \mathbb{E}[(\mathfrak{Z}_i^{t-1})^2] + \mathbb{P}[\mathfrak{Z}_i^{t-1} \neq \mathfrak{Z}_j^{t-1}] \mathbb{E}[\mathfrak{Z}_i^{t-1} \mathfrak{Z}_j^{t-1}]) \\ &= \left(\frac{t-1}{t}\right)^2 \left(\frac{1}{n} \frac{1}{t-1} + \frac{n-1}{n} \mathbb{E}[Z_i^{t-1} Z_j^{t-1}]\right). \end{aligned}$$

Now let $u_t = t^2 \mathbb{E}[Z_i^t Z_j^t] - t + n$. The last line becomes $u_t = \frac{n-1}{n} u_{t-1}$. Since $u_1 = n - 1$, we have

$$\begin{aligned} \mathbb{E}[Z_i^t Z_j^t] &= \frac{u_t + t - n}{t^2} = \frac{\left(\frac{n-1}{n}\right)^{t-1} (n-1) + t - n}{t^2} \\ &= \frac{1}{t^2 n^{t-1}} \{(n-1)^t - n^t + t n^{t-1}\} = \frac{t-1}{2t} \frac{n^{t-2}}{n^{t-1}} + \dots = O_n(n^{-1}). \end{aligned}$$

In conclusion,

$$\begin{aligned} \text{var}(\bar{Z}_i^t) &= \frac{1}{nt} + \frac{n(n-1)}{n^2} \frac{1}{t^2 n^{t-1}} \{(n-1)^t - n^t + t n^{t-1}\} \\ &= \frac{1}{nt} \left[1 + \frac{n(n-1)}{t} \{(1 - 1/n)^t - 1 + t/n\} \right]. \end{aligned}$$

For n fixed, and $t \rightarrow +\infty$, $t \text{var}(\bar{Z}_i^t) \rightarrow 1$, a limit that does not depend on n , i.e. the system eventually degenerates to a single path. If we set $n = ct$, then $n \text{var}(\bar{Z}_i^t) \rightarrow C$, for some $C > 0$. Bearing in mind that the actual posterior variance should be $O(t^{-1})$, this means that, to bound the *relative error* uniformly over a given time interval, i.e. for $t = 1, \dots, T$, one must take $n = O(T)$.

6 On the degeneracy of path functionals in SMC (Chopin and Schäfer)

Much of the confusion around the degeneracy of particle learning and similar algorithms (Fearnhead, 2002, Storvik, 2002) seems related to the lack of formal results regarding the degeneracy of path functional in Sequential Monte Carlo. We'd like to report here some preliminary investigation on this subject.

Consider a standard state-space model, with observed process (y_t) , and hidden Markov process (x_t) , and a basic particle filter, which would track the complete trajectory $x_{1:t}$, i.e. which would produce, at each iteration t , N simulated trajectories $x_{1:t}^{(n)}$, with some weight $w_t^{(n)}$, so as to approximate $p(x_{1:t}|y_{1:t})$. It is well known that the Monte Carlo error regarding the expectation of $\varphi(x_{1:t})$ (a) remains bounded over time if $\varphi(x_{1:t}) = x_t$, (the filtering problem), and (b) blows away, at an exponential rate, if $\varphi(x_{1:t}) = x_1$ (the smoothing problem). Chopin (2004) formalises these two statements by studying the asymptotic variance that appears in the central limit theorem for the corresponding particle estimates.

As mentioned above, and to the best of our knowledge, there is currently no formal result on the divergence of the asymptotic variance for test functions like $\varphi(x_{1:t}) = t^{-1} \sum_{i=1}^t x_i$, i.e. some symmetric function with respect to the complete trajectory. (The fact that this function is a sufficient statistic should not play any role in this convergence study.) One difficulty is that the iterative definition of the asymptotic variance given by Chopin (2004) leads to cumbersome calculations.

We managed however to compute this asymptotic variance exactly, for the Gaussian local level model:

$$x_{t+1} | x_t \sim N(x_t, 1), \quad y_t | x_t \sim N(x_t, 1)$$

and the functional $\varphi(x_{1:t}) = t^{-1} \sum_{i=1}^t x_i$. In this case, the asymptotic variance diverges at rate $O(e^{ct}/t^2)$. Exact calculations may be requested from the authors. We plan to extend these results to a slightly more general model, e.g. with unknown variances, and a function φ which would be a sufficient statistic for such parameters. We conjecture that this exponential divergence occurs for many models: basically, in an average like $\varphi(x_{1:t}) = t^{-1} \sum_{i=1}^t x_i$, the Monte Carlo error attached to x_1/t should be $O(e^{ct}/t^2)$, and should dominate all the other terms. This is at least what one observes in toy examples. After, say, 100 iterations of a particle filter, the number of distinct values within all the simulated trajectories (that have survived so far) for the component x_1 is typically very small, and the degeneracy in the x_1 dimension seems sufficient to endanger the accuracy of any estimate based on the complete trajectory $x_{1:t}$.

7 Remarks on the rejoinder (Robert)

Lopes et al. (2010) published a rejoinder on the discussions of their paper and the following is a detailed examination of the arguments found in this rejoinder, which requires a preliminary reading of the above papers as well as our discussion. (All quotes are taken verbatim from the rejoinder.)

“Particle learning based on the product estimate and MCMC based on Chib’s formula produce relatively similar results either for small or large samples.”

This statement about the estimation of the marginal likelihood (or the evidence) and the example A that is associated with it in the rejoinder thus comes to contradict our (rather intensive) simulation experiment which, as reported in the discussion (Section 2), concludes to a strong bias in evidence approximation induced by using particle learning, whether or not the product estimator is used. We observed there that there were two levels of degeneracy, one due to the product solution (errors in a product being more prone to go and multiply) and one due to the particle nature of the sequential method (which does not refresh particles from earlier periods). Figures 2–5 are at odds with the one presented in the rejoinder, maybe because we consider 10,000 observations rather than 100. (I also fail to understand how the “Log-predictive (TRUE)” quantity is derived.)

“Black-box sequential importance sampling algorithms and related central limit theorems are of little use in practice.”

This is a quote from the rejoinder that is rather puzzling. There is nothing wrong with the central limit theorem which is the basis of error assessment in Monte Carlo studies (Robert and Casella, 2009). Indeed, one major consequence of the central limit theorem is that it provides a precise scale for the speed of convergence of Monte Carlo estimates and thus an indicator on the number of particles needed for a given precision level. The authors of the rejoinder then criticise our use of “1000 particles in 5000 dimensional problems” as we “shouldn’t be surprised at all with some of our findings”. This is not factually exact since I find no trace in the discussion of such a case: we use 10,000 particles in all examples and the target is either the distribution of the 4 mixture parameters, the evidence or the distribution of a one-dimensional sufficient statistic. Furthermore, these values of n and N are those used in their example D. More importantly, nor the paper neither the rejoinder map a practical strategy on how to increase the computational effort along with the number of observations.

“This argument [that the Monte Carlo variance will ‘blow-up’] is incorrect and extremely misleading.”

This point is central both to the discussions above and to the rejoinder, as the authors maintain that the inevitable particle degeneracy does not impact the distribution of the sufficient statistics. The argument about using time averages over particle paths rather than sums appears reasonable at first. Actually, taking an empirical average in almost stationary situations should produce an approximately normal distribution. With an asymptotic variance different from 0 (thanks to the central limit theorem) However, this is not the main argument used in the discussions. Degeneracy in the particle path means that the early terms in the average are less and less diverse in the sample average. Therefore it is not that surprising that the variance is decreasing down to too small a value! As shown in Figure 8 above, degeneracy due to resampling may induce severe biases in the distribution of empirical averages while giving the impression of less variability (which is a recurrent argument in the rejoinder). Furthermore, the fact that parameters are simulated [rather than fixed] in the particle filter means that the process is not geometrically ergodic, hence that Monte Carlo errors tend to accumulate along iterations, rather than compensate. (This is why the comparison between PL and sampling importance resampling is particularly relevant, because it does not address this accumulation.) The rejoinder also quotes Olsson et al. (2008) for justifying the decrease in the Monte Carlo variance. This is somehow surprising in that (a) Olsson et al. (2008) show that there is degeneracy without a fixed-lag smoothing and (b) they require a geometric forgetting property on the filtering dynamics. In addition, I think that Example E used to illustrate the point about variance reduction is not very appropriate for this issue because the hidden Markov chain is a Gaussian random walk, hence cannot be stationary (a fact noted by the authors). And once again a decrease in the “MC error” does not mean a converging algorithm because degeneracy naturally induces empirical variance decrease. (I also fail to see why the “prior” on (x_t) is improper.) The final (if recurrent) argument that “PL parameters do not degenerate” is somehow puzzling: by nature, those parameters are simulated from a distribution conditional on the sufficient parameters. So obviously the simulated parameters all differ. But this does not mean that they are marginally distributed from the right distribution.

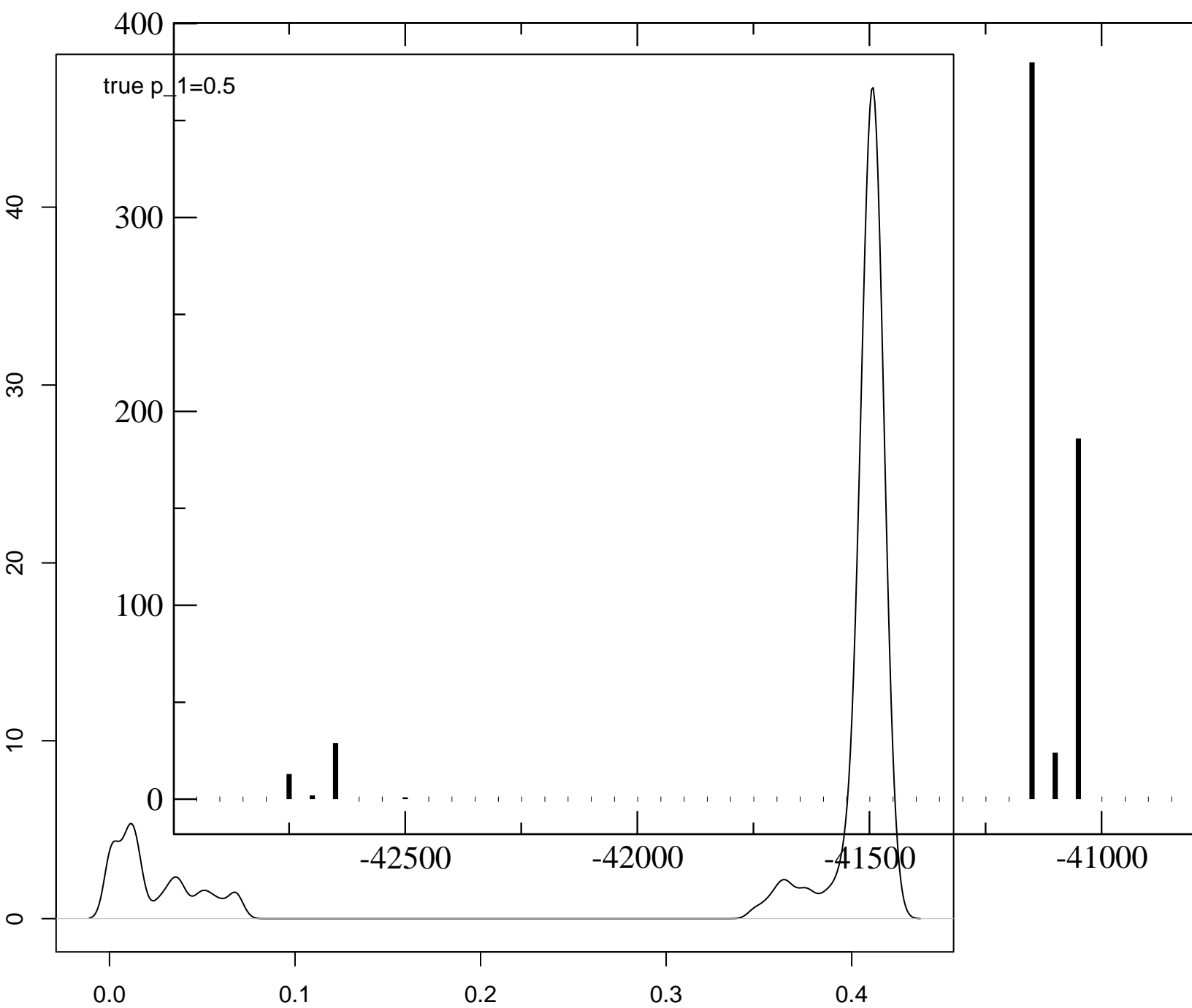
“MCMC schemes depend upon the not so trivial task of assessing convergence. How long should the burn-in G_0 be?”

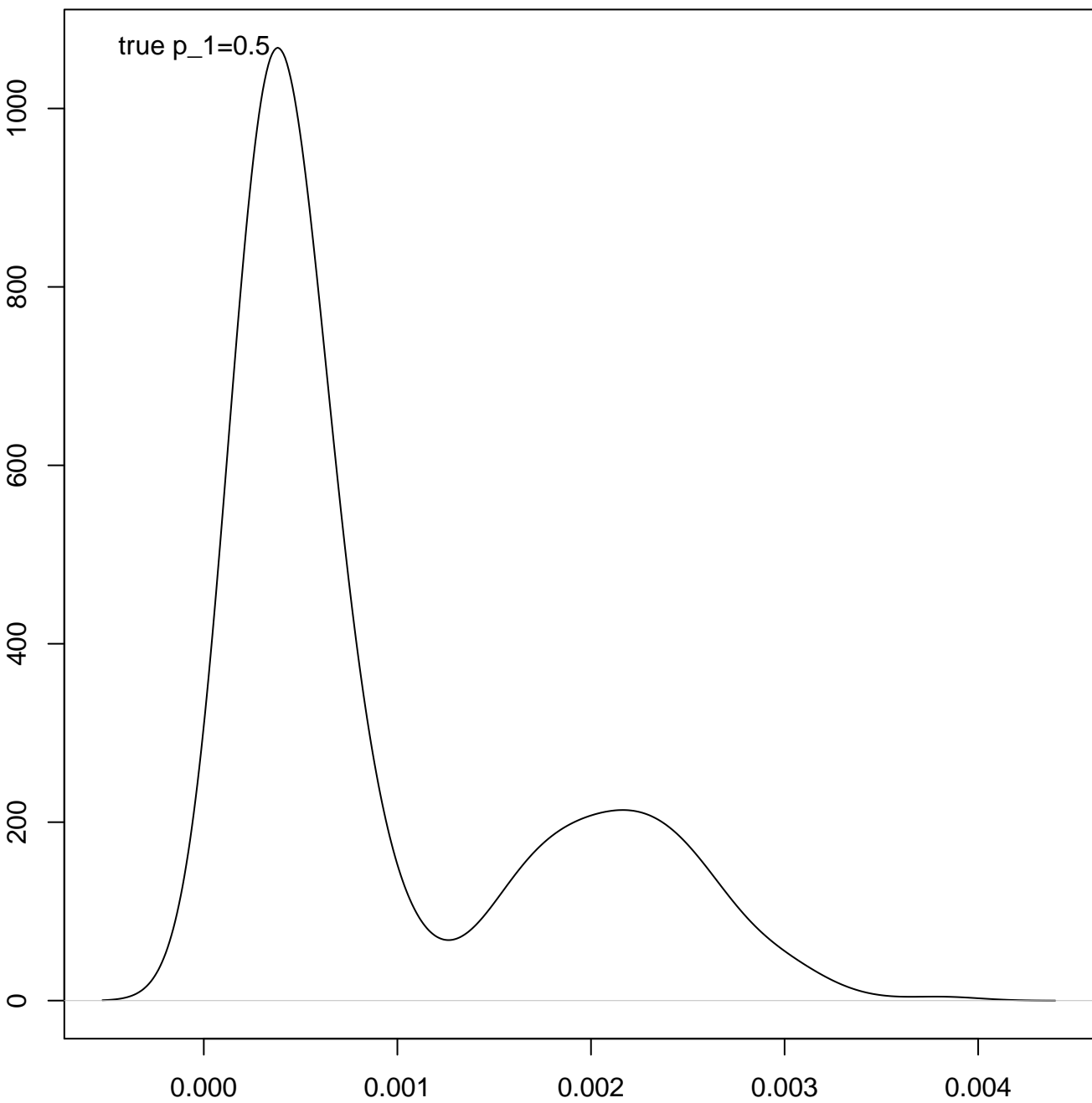
The rejoinder concludes with recommendations that sound more like a drafted to-do note the authors forgot to remove than an accumulation of true recommendations. (The above quote rather clearly supports our first point in the discussion.) It seems to me that the comparison between MCMC and particle filters is not particularly relevant, simply because particle filters apply in [sequential] settings where MCMC cannot be implemented. To try to promote PL over MCM by arguing that MCMC produces dependent draws while having convergence troubles is not needed (besides, PL also produces [unconditional] dependent draws). To advance that the Monte Carlo error for PL is of order C_T/\sqrt{N} is not relevant either because C_T is exponential in T and because MCMC also has an error in \sqrt{N} .

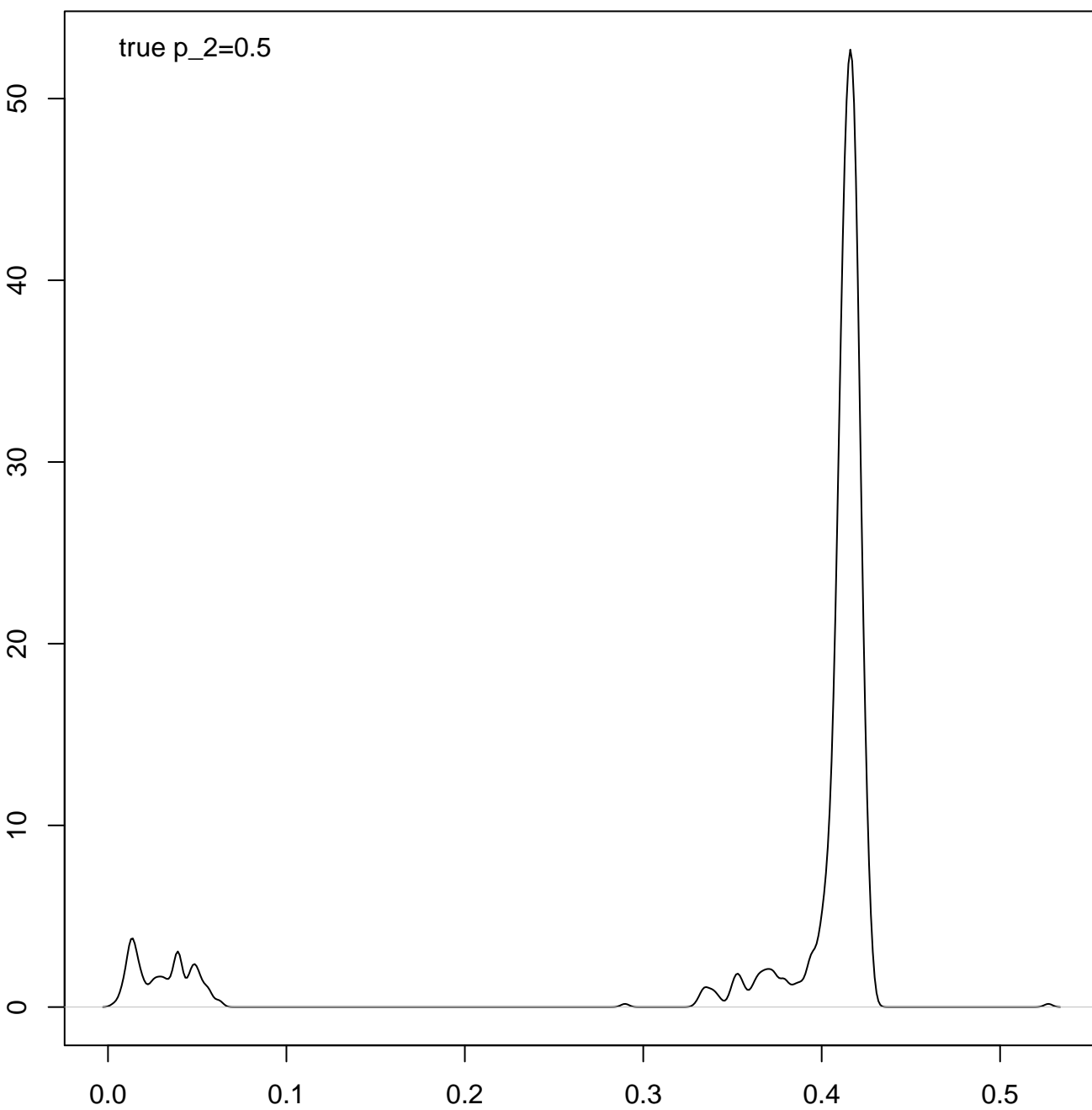
References

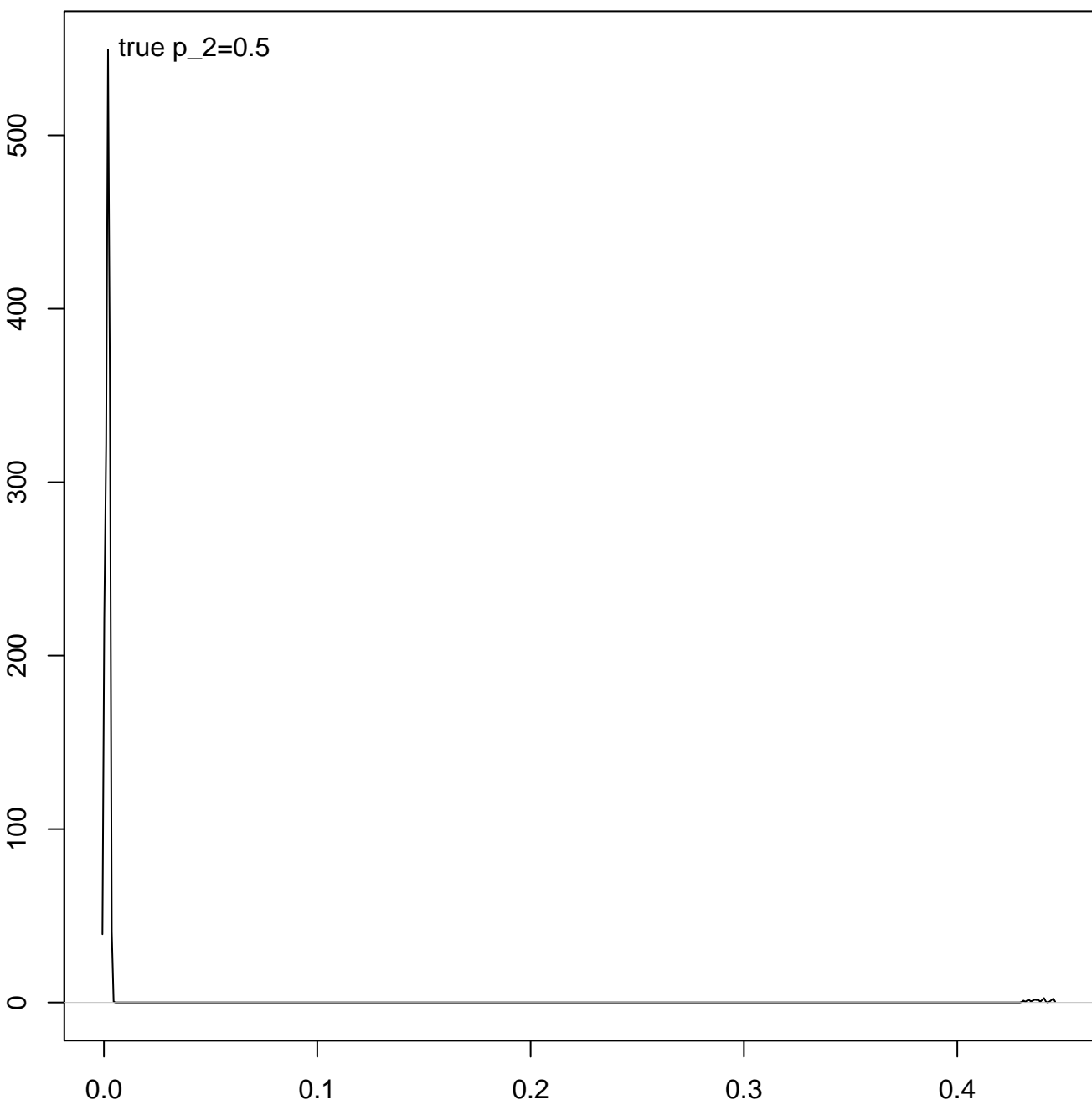
- Berkhof, J., I. van Mechelen, and A. Gelman. 2003. A Bayesian approach to the selection and testing of mixture models. *Statistica Sinica* 13: 423–442.
- Cappé, O., E. Moulines, and T. Rydén. 2004. *Hidden Markov Models*. Springer-Verlag, New York.
- Carvalho, C., H. Lopes, N. Polson, and M. Taddy. 2009. Particle learning for general mixtures. Tech. Rep. 09-02, Duke University.
- Chen, M., Q. Shao, and J. Ibrahim. 2000. *Monte Carlo Methods in Bayesian Computation*. Springer-Verlag, New York.
- Chib, S. 1995. Marginal likelihood from the Gibbs output. *J. American Statist. Assoc.* 90: 1313–1321.
- Chopin, N. 2004. Central Limit theorem for sequential Monte Carlo methods and its application to Bayesian inference. *Ann. Statist.* 32(6): 2385–2411.
- Chopin, N. and C. Robert. 2010. Properties of nested sampling. *Biometrika* To appear, doi:10.1093/biomet/asq021.
- Darmois, G. 1935. Sur les lois de probabilité à estimation exhaustive. *Comptes Rendus Acad. Sciences Paris* 200: 1265–1266.
- Del Moral, P., A. Doucet, and A. Jasra. 2006. Sequential Monte Carlo samplers. *J. Royal Statist. Society Series B* 68(3): 411–436.
- Diebolt, J. and C. Robert. 1990. Estimation des paramètres d’un mélange par échantillonnage bayésien. *Notes aux Comptes-Rendus de l’Académie des Sciences I* 311: 653–658.
- . 1994. Estimation of finite mixture distributions by Bayesian sampling. *J. Royal Statist. Society Series B* 56: 363–375.
- Douc, R., O. Cappé, E. Moulines, and C. Robert. 2002. On the convergence of the Monte Carlo maximum likelihood method for latent variable models. *Scandinavian J. Statist.* 29(4): 615–636.
- Fearnhead, P. 2002. MCMC, sufficient statistics and particle filters. *J. Comp. Graphical Statist.* 11: 848–862.
- Gordon, N., J. Salmond, and A. Smith. 1993. A novel approach to non-linear/non-Gaussian Bayesian state estimation. *IEEE Proceedings on Radar and Signal Processing* 140: 107–113.
- Kitagawa, G. 1996. Monte Carlo filter and smoother for non-Gaussian non-linear state space models. *J. Comput. Graph. Statist.* 5: 1–25.
- Lopes, H., C. Carvalho, M. Johannes, and N. Polson. 2010. Particle learning for sequential Bayesian computation (with discussion and rejoinder). In *Bayesian Statistics 9*, eds. J. Bernardo, J. B. M.J. Bayarri, A. Dawid, D. Heckerman, A. Smith, and M. West. Oxford University Press. To appear.

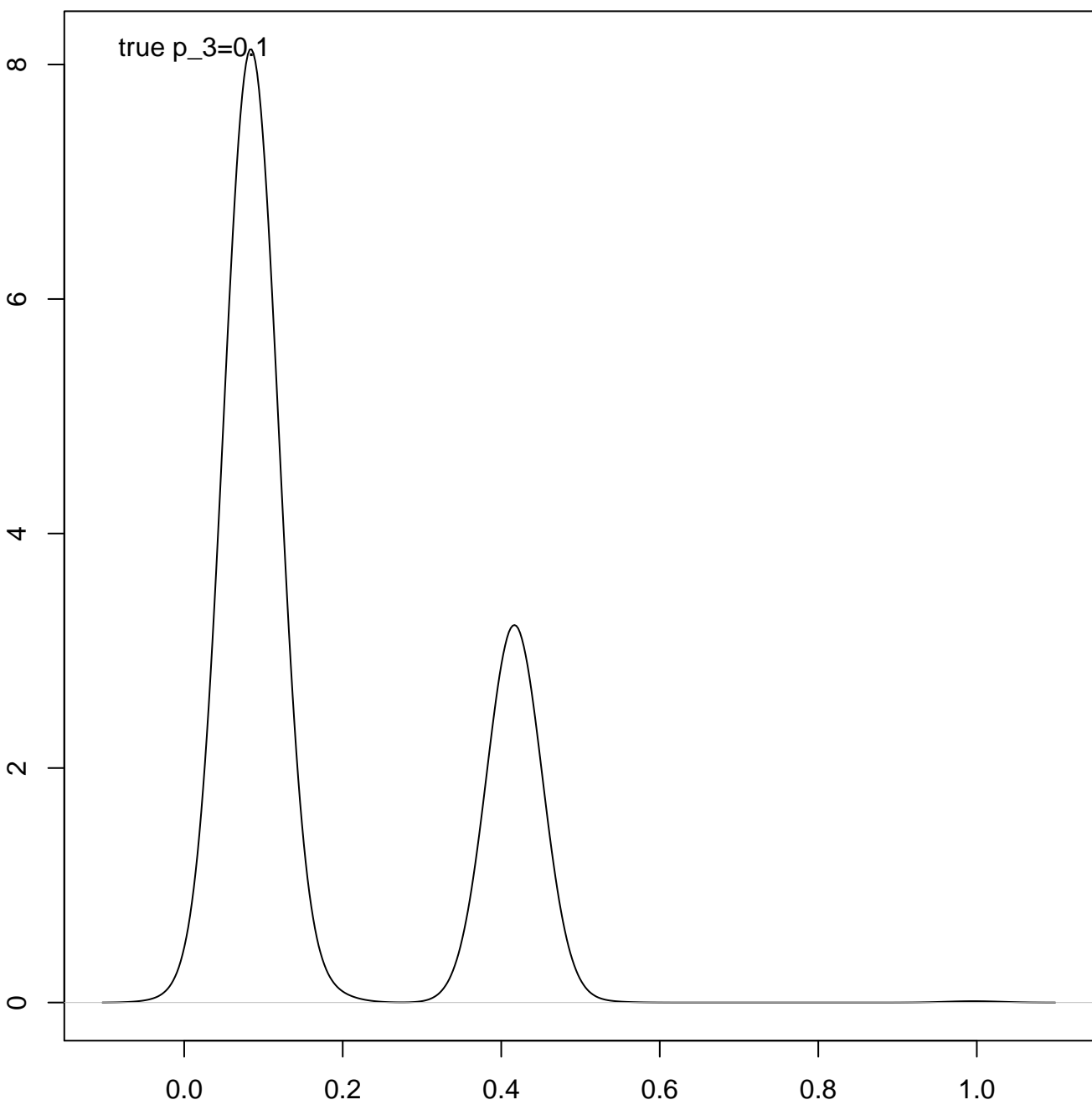
- Marin, J. and C. Robert. 2010. Importance sampling methods for Bayesian discrimination between embedded models. In *Frontiers of Statistical Decision Making and Bayesian Analysis*, eds. M.-H. Chen, D. Dey, P. Müller, D. Sun, and K. Ye. Springer-Verlag, New York. To appear.
- Olsson, J., O. Cappé, R. Douc, and E. Moulines. 2008. Sequential Monte Carlo smoothing with application to parameter estimation in nonlinear state space models. *Bernoulli* 14(1): 155–179.
- Pitt, M. and N. Shephard. 1999. Filtering via simulation: auxiliary particle filters. *J. American Statist. Assoc.* 94(446): 590–599.
- Robert, C. and G. Casella. 2009. *Introducing Monte Carlo Methods with R*. Springer-Verlag, New York.
- Rubin, D. 1988. Using the SIR algorithm to simulate posterior distributions. In *Bayesian Statistics 3: Proceedings of the Third Valencia International Meeting, June 1-5, 1987*, eds. J. Bernardo, M. Degroot, D. Lindley, and A. Smith. Clarendon Press.
- Storvik, G. 2002. Particle filters for state space models with the presence of static parameters. *IEEE Trans. Signal Process.* 50: 281–289.

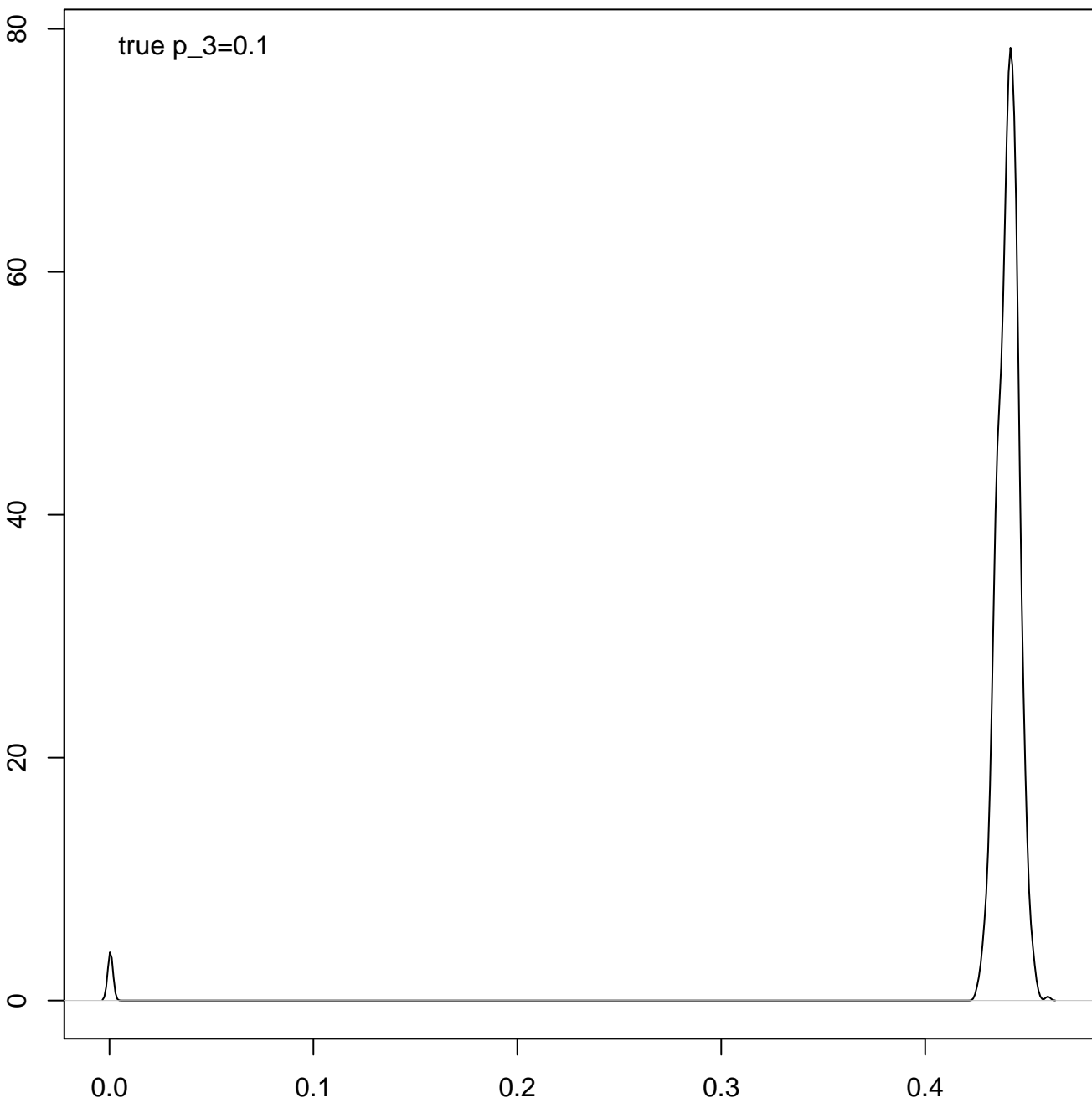


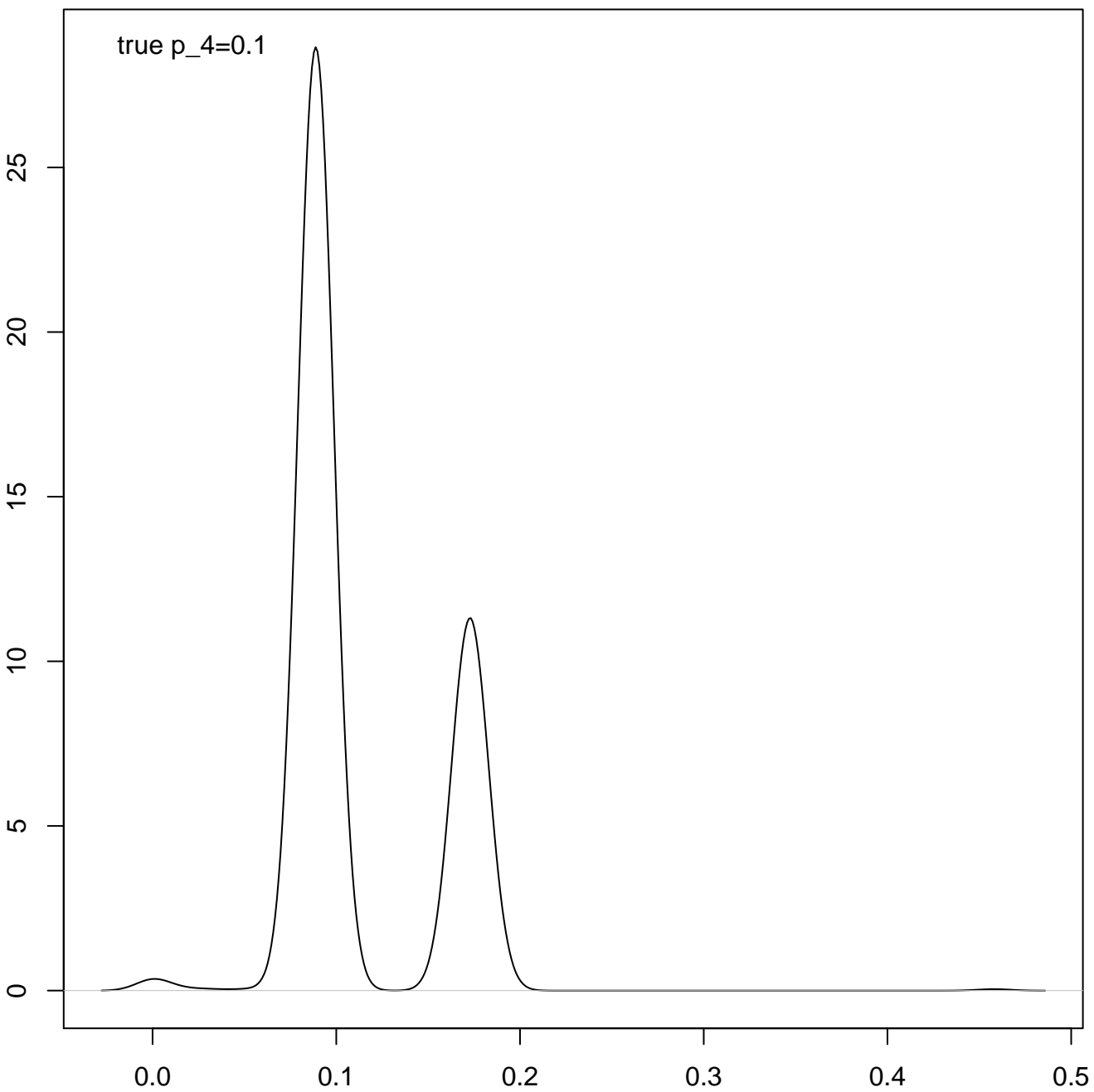


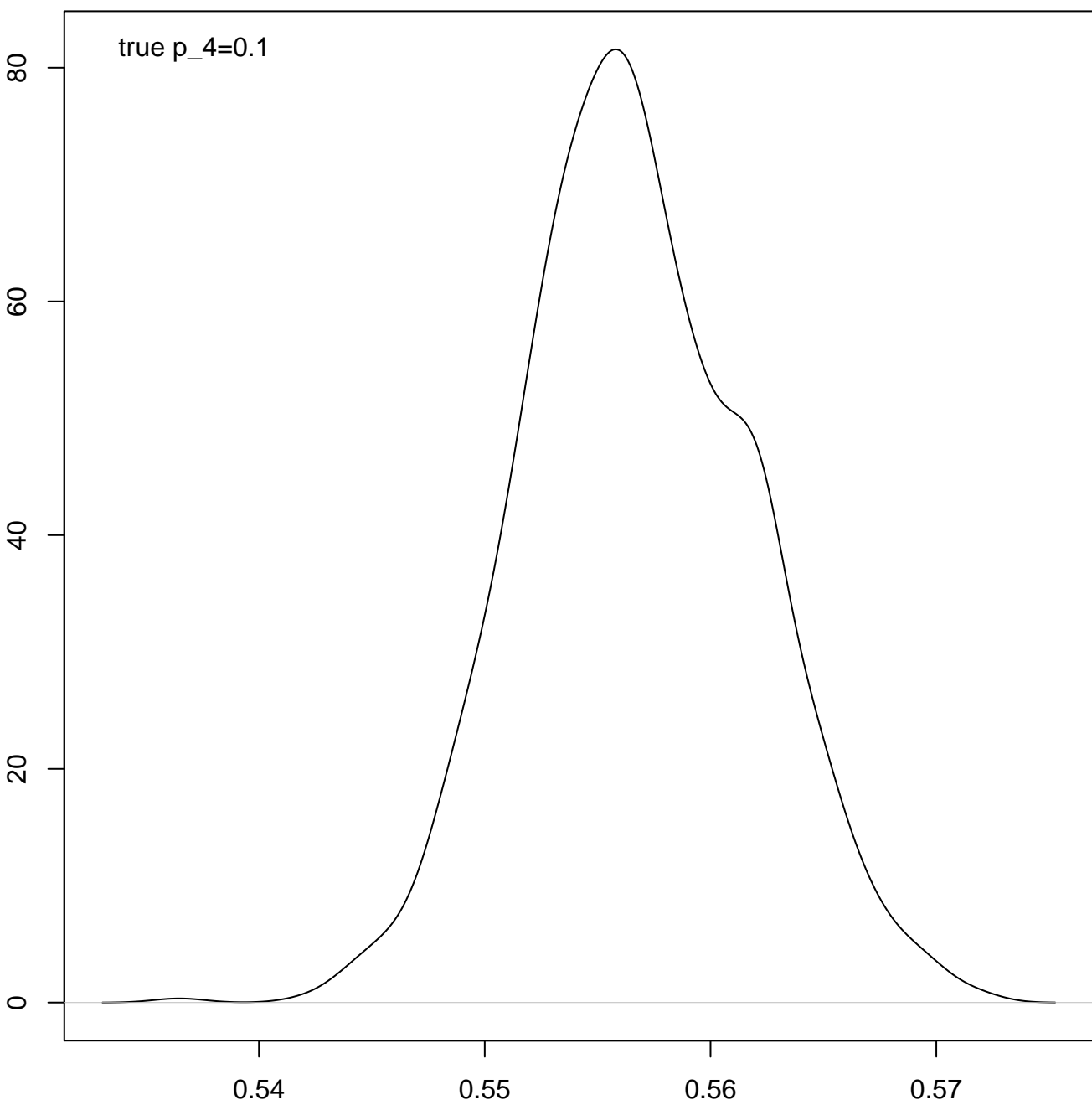


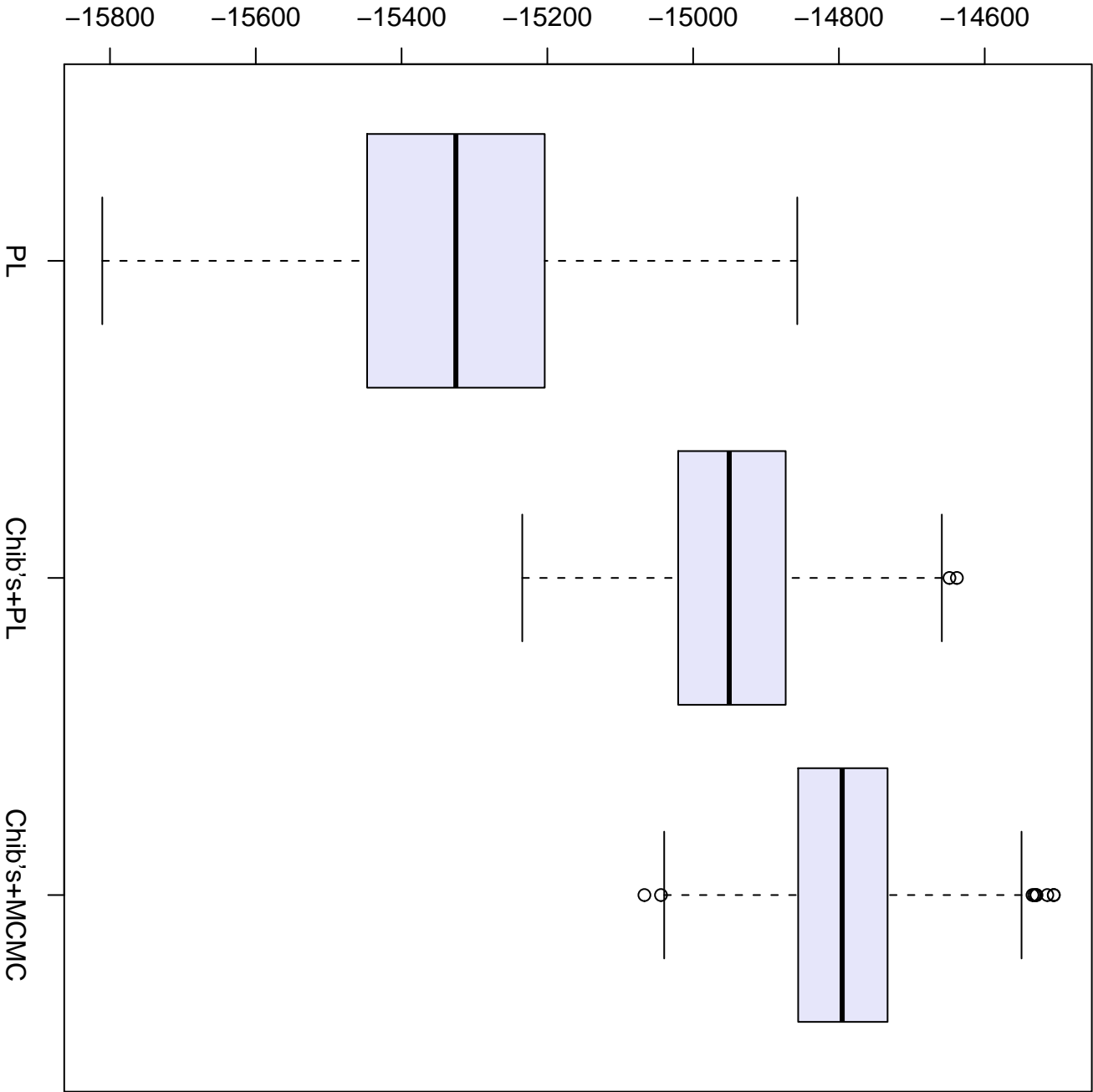


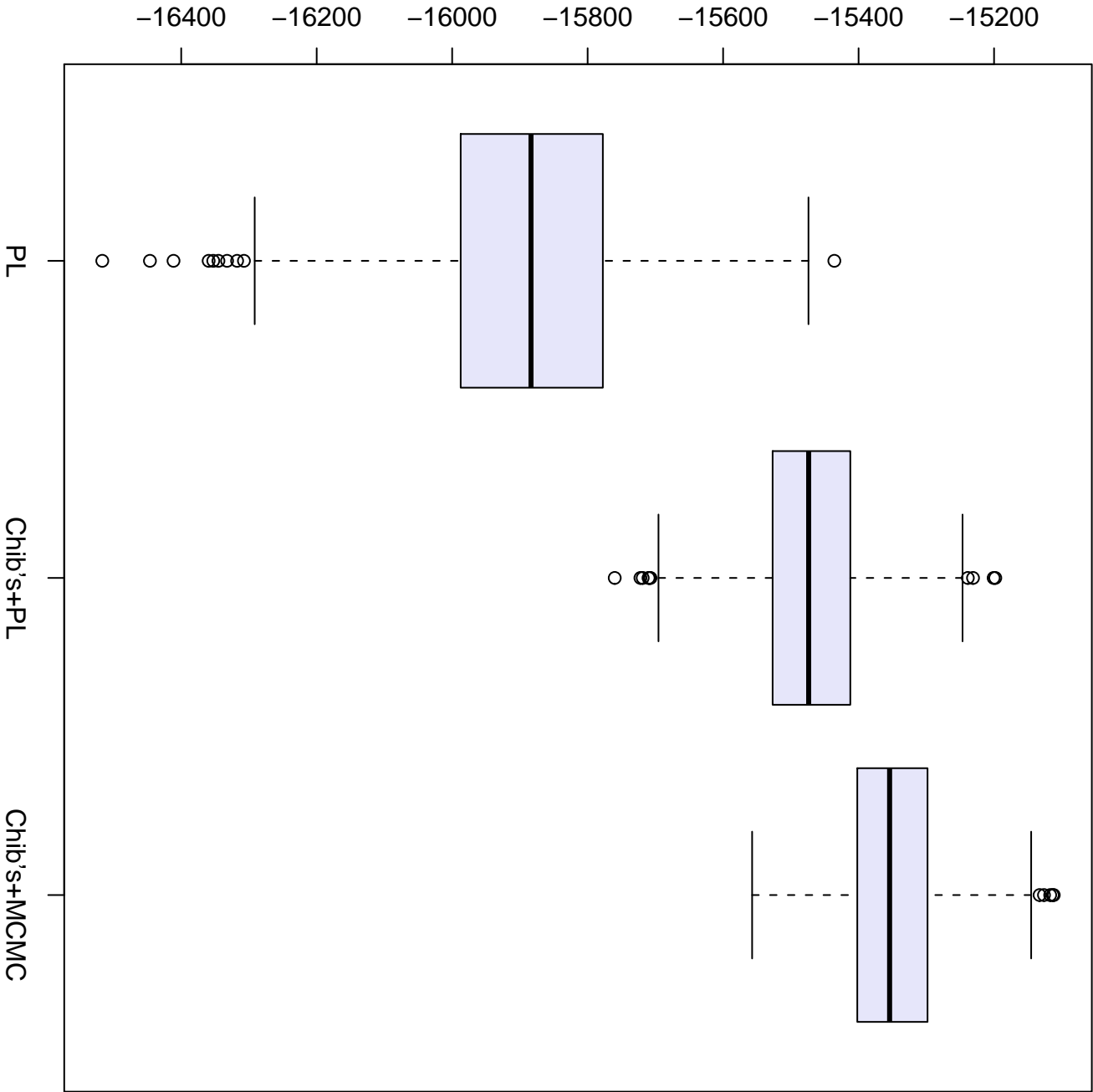


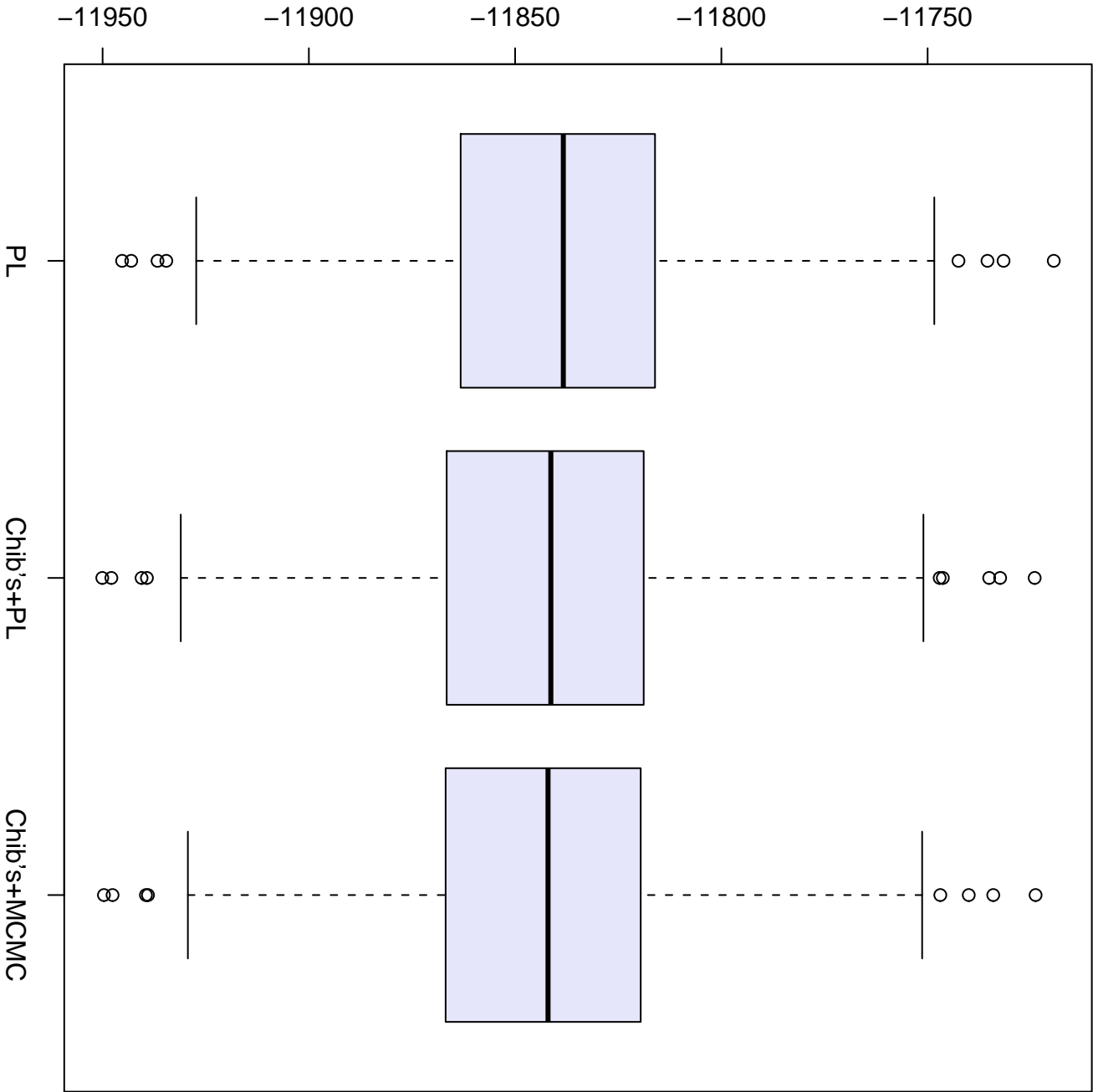


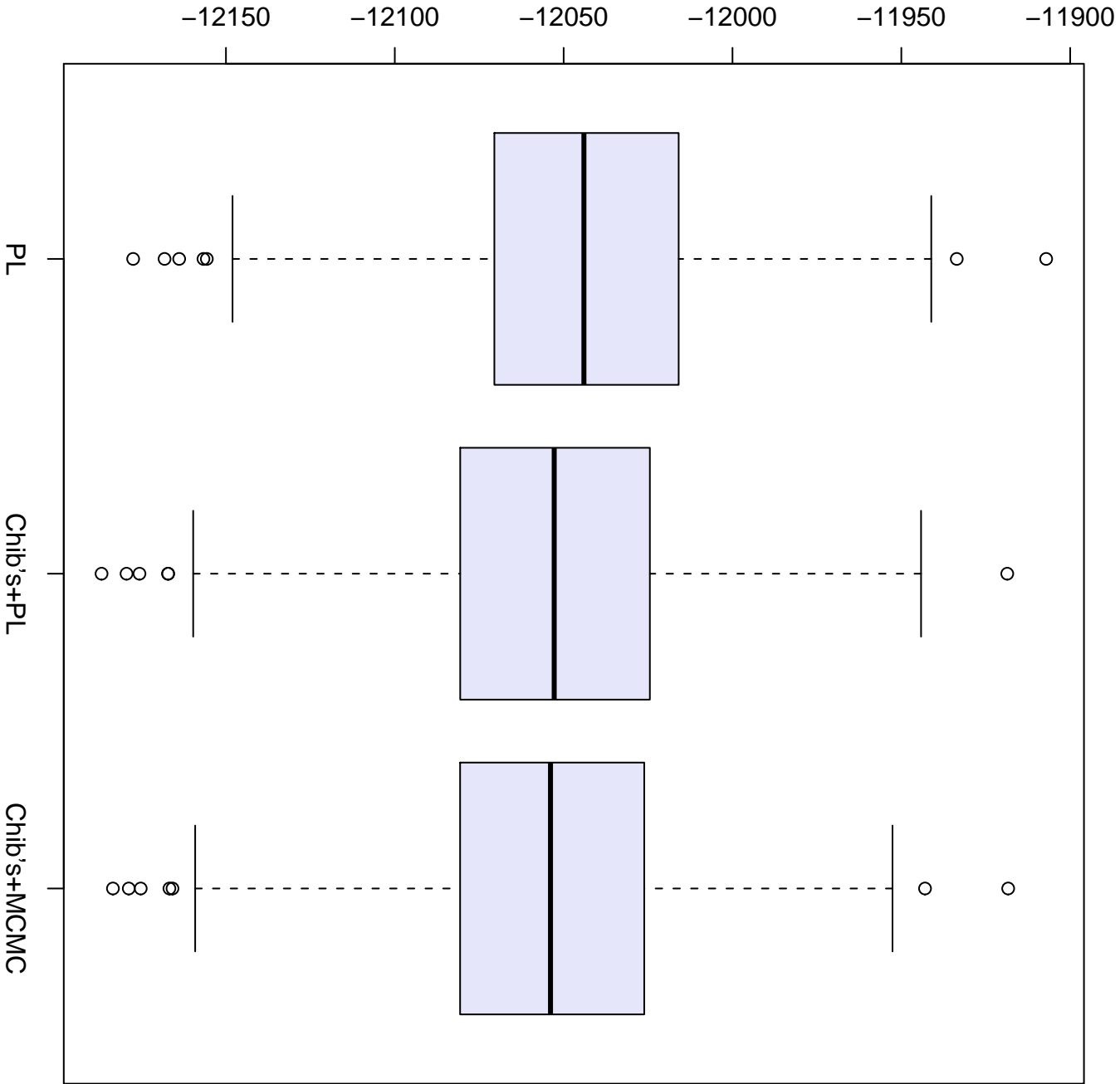


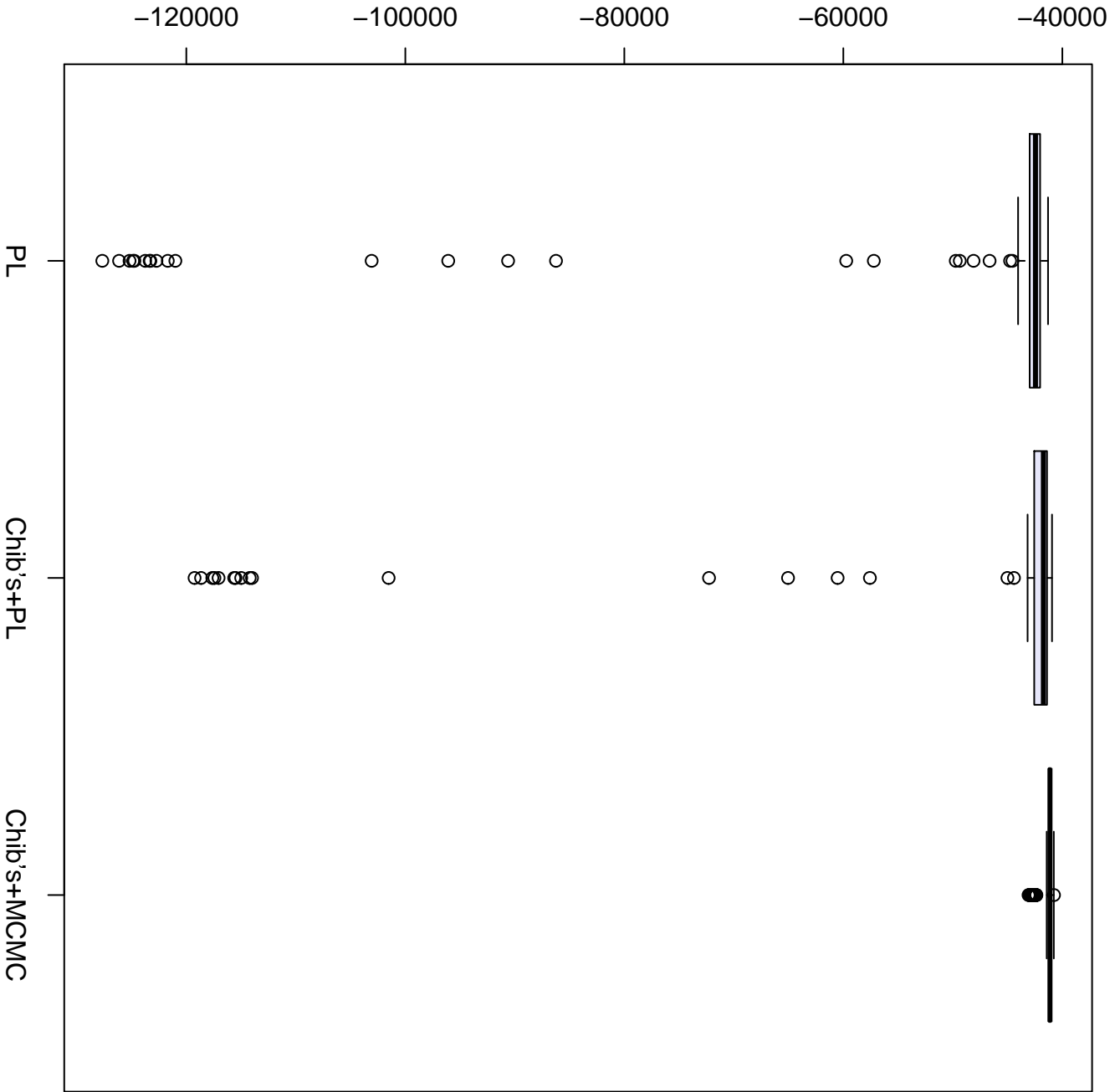


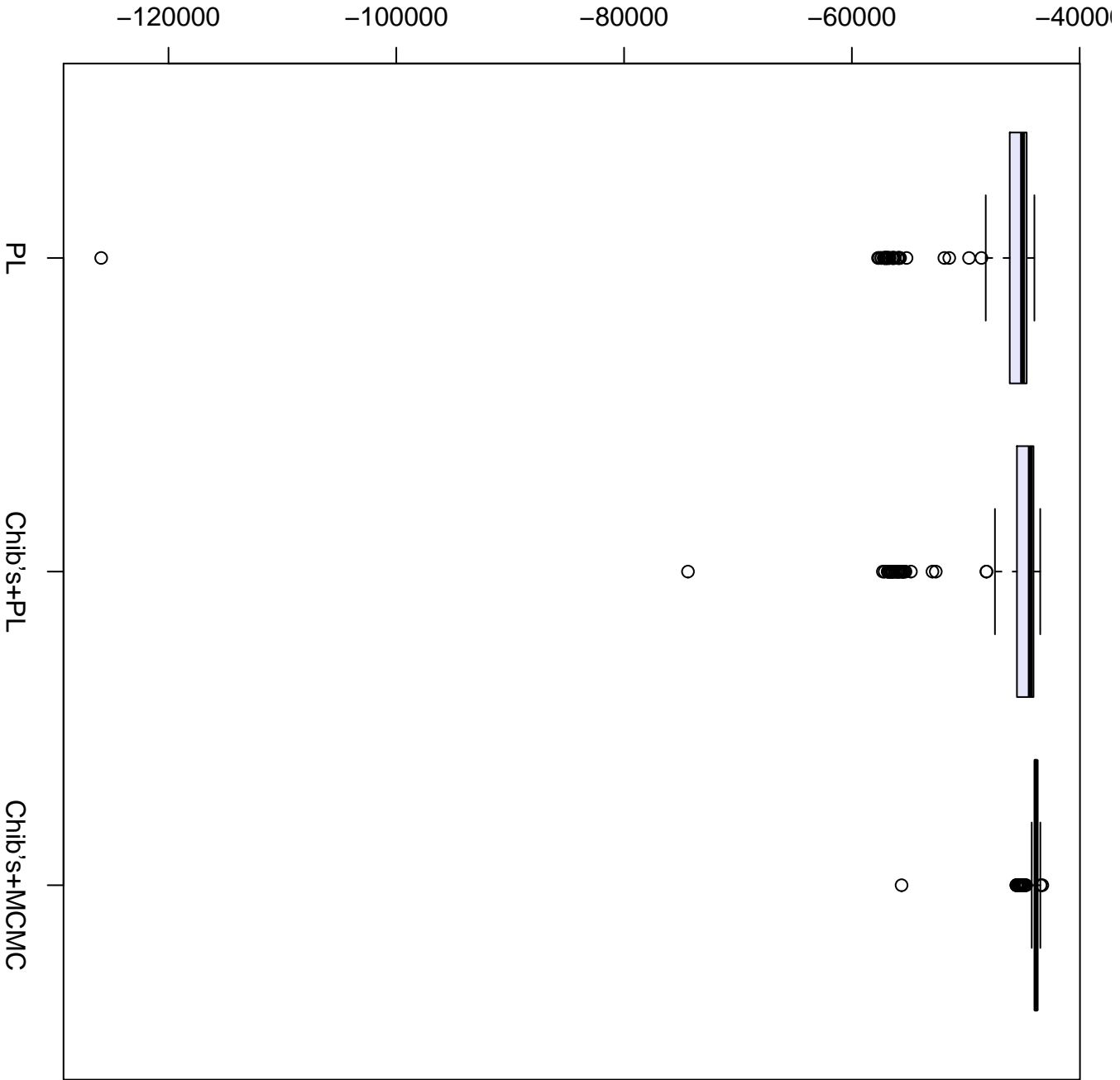




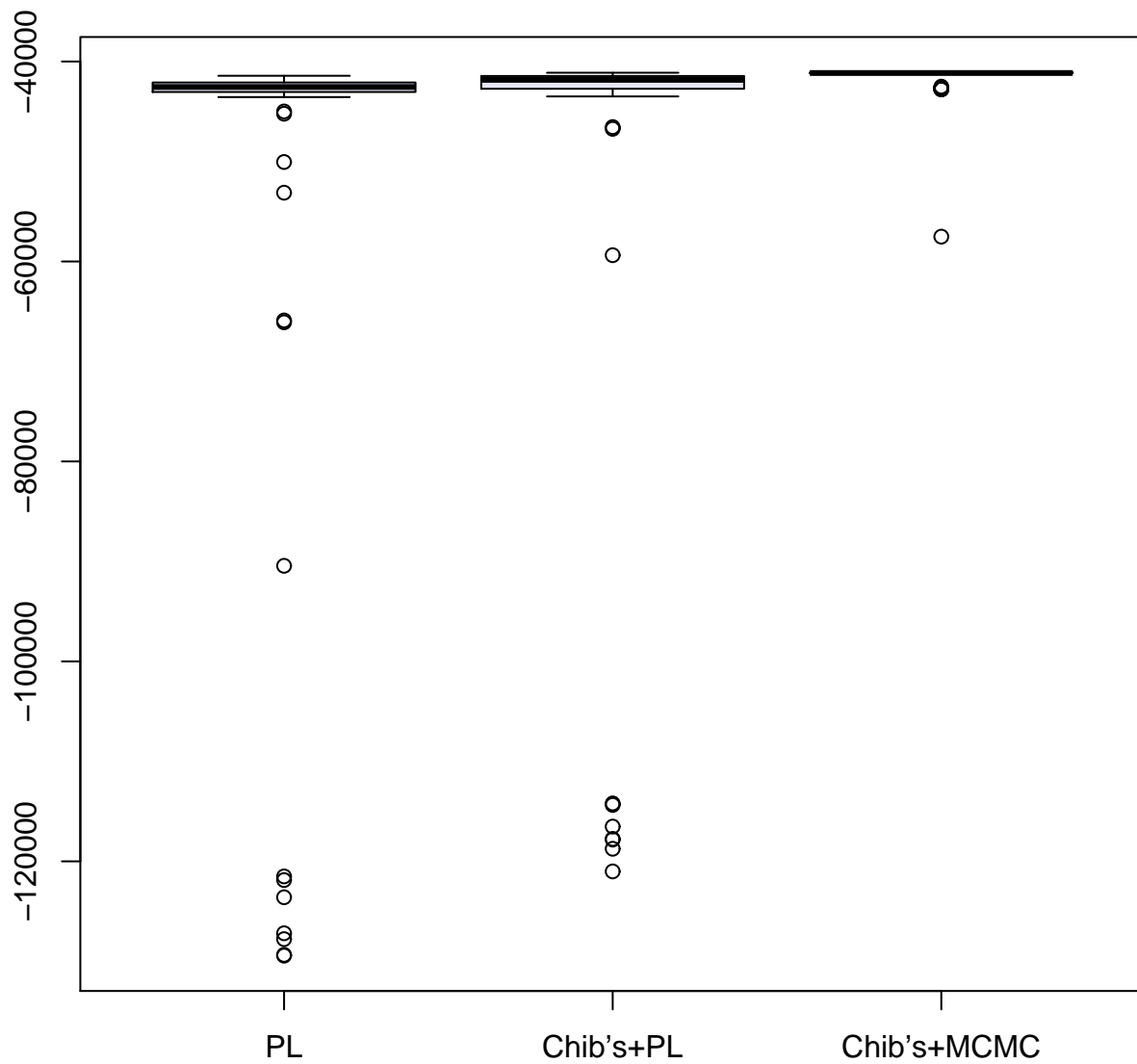




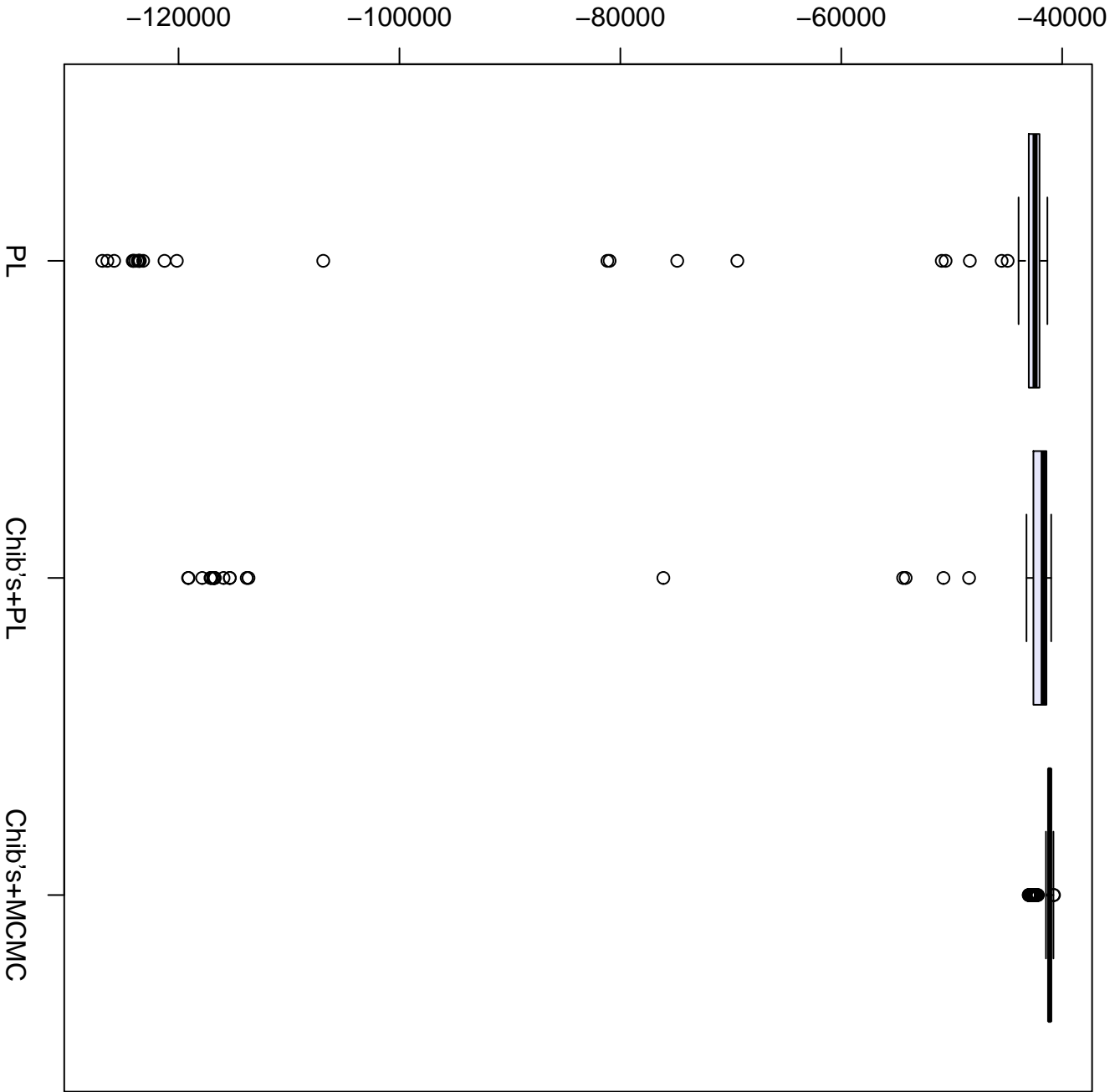


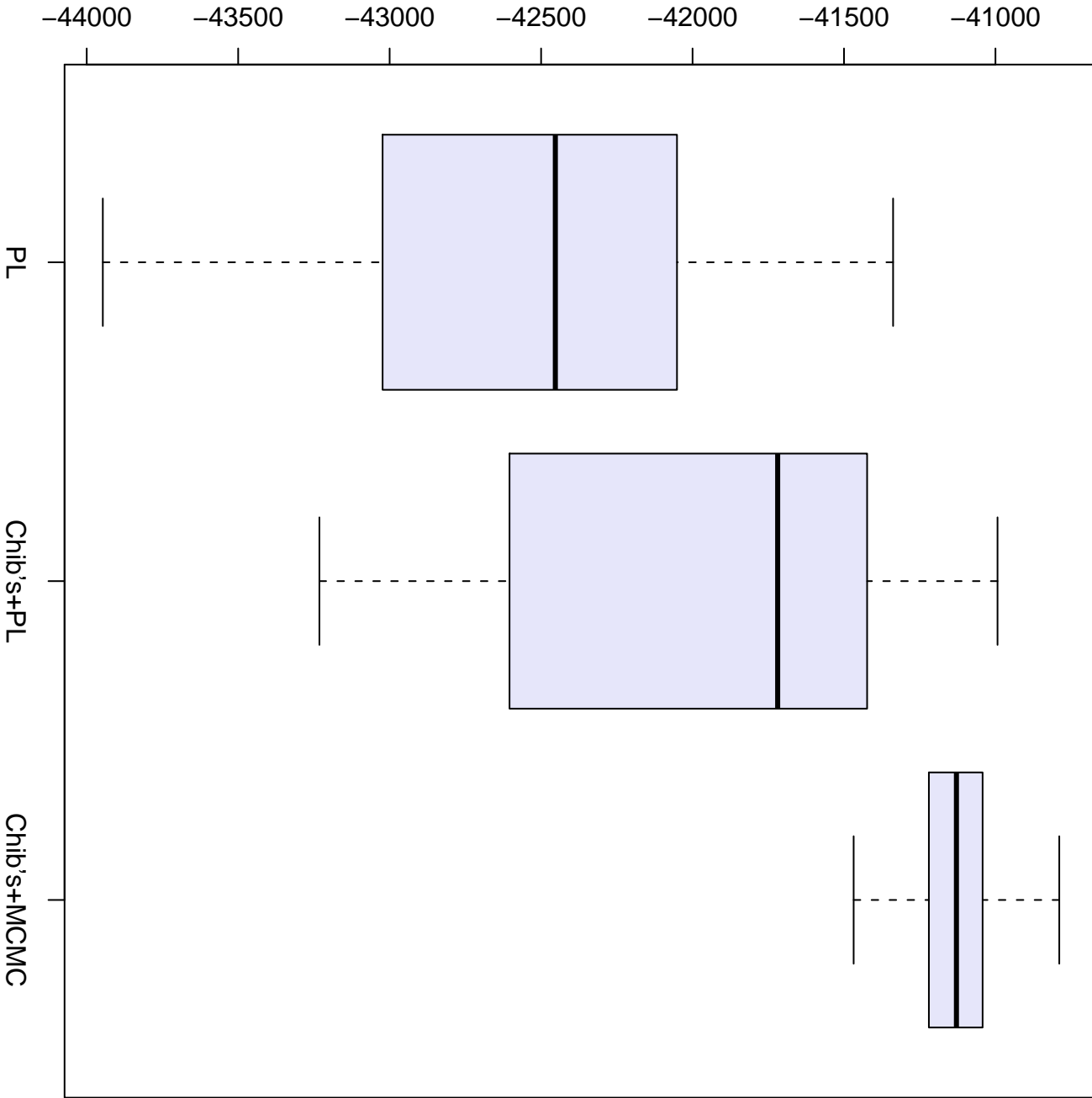


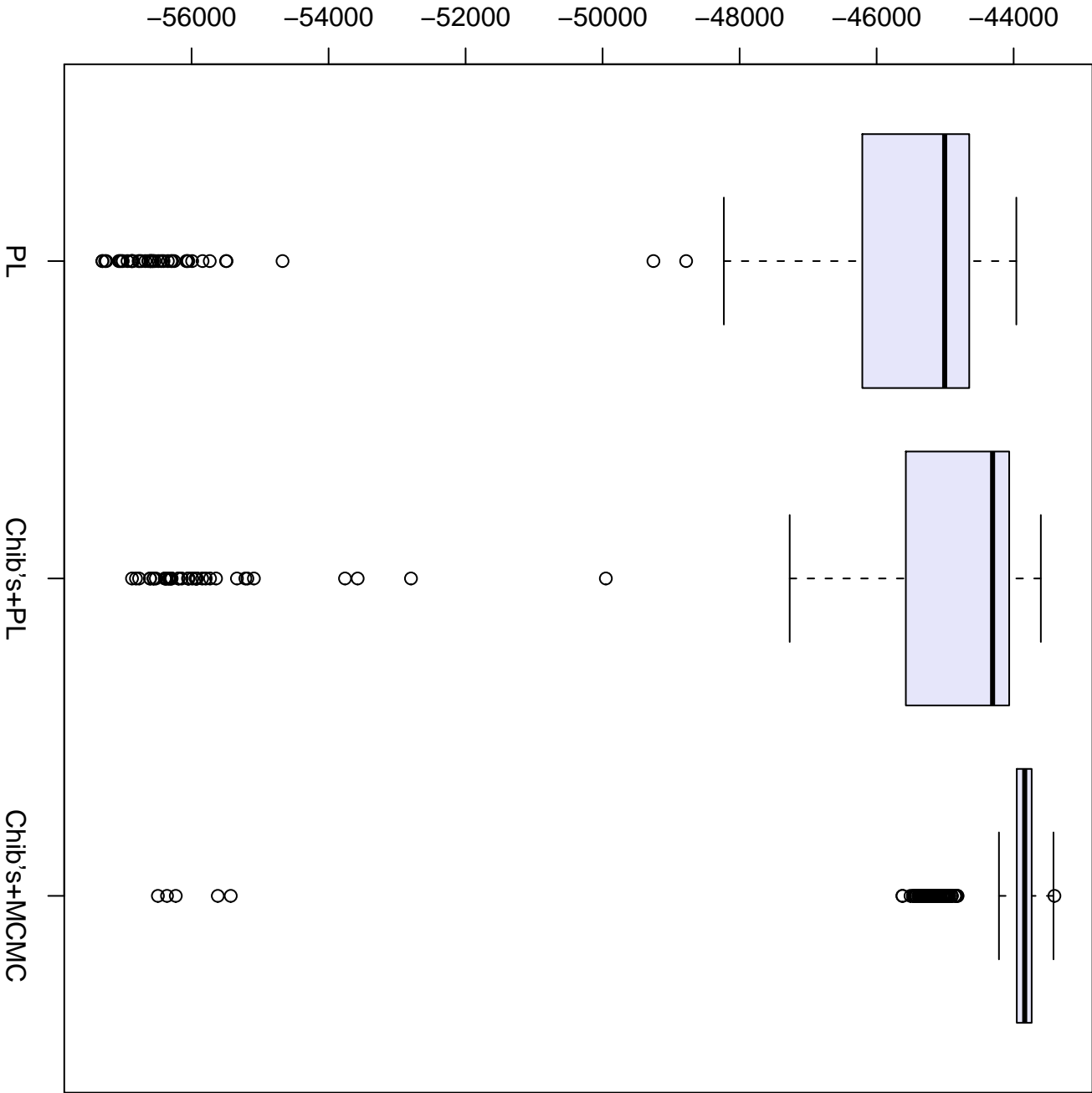
same data; $\lambda=(10,30,110,150)$

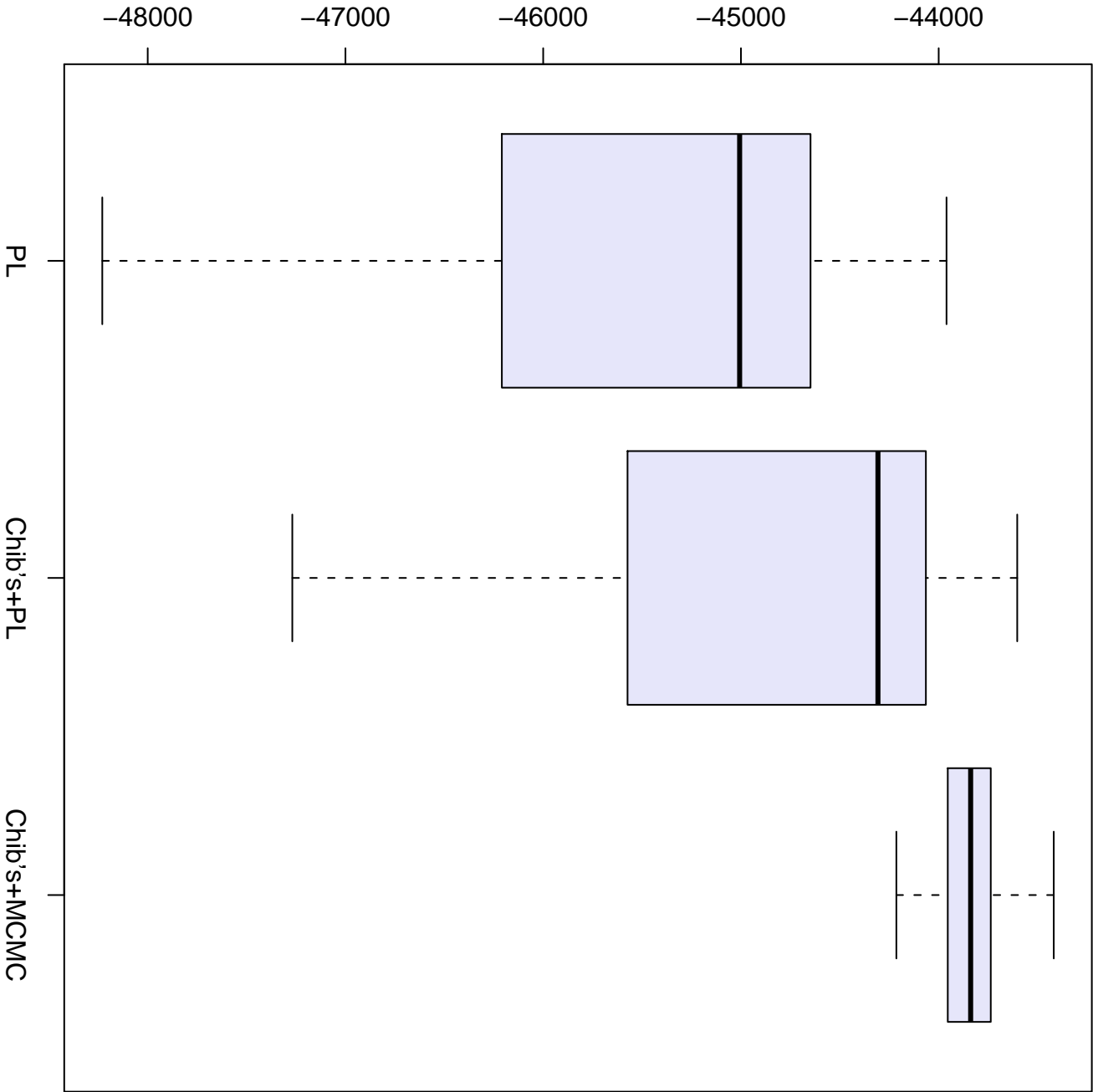


nobs= 10000









-34000 -33500 -33000 -32500 -32000 -31500 -31000

PL

Chib's+PL

Chib's+MCMC

